

# Stochastic gradient updates yield deep equilibrium kernels

Anonymous authors

Paper under double-blind review

## Abstract

Implicit deep learning allows one to compute with implicitly defined features, for example features that solve optimisation problems. We consider the problem of computing with implicitly defined features in a kernel regime. We call such a kernel a deep equilibrium kernel (DEK). Specialising on a stochastic gradient descent (SGD) update rule applied to features in a latent variable model, we find an exact corresponding deterministic update rule for the DEK in a high dimensional limit. This derived update rule resembles previously introduced infinitely wide neural network kernels. To perform our analysis, we describe an alternative parameterisation of the link function of exponential families, a result that may be of independent interest. This new parameterisation allows us to draw new connections between a statistician’s inverse link function and a machine learner’s activation function. We describe an interesting property of SGD in this high dimensional limit: even though individual iterates are random vectors, inner products of any two iterates are deterministic, and can converge to a unique fixed point as the number of iterates increases. We find that the DEK empirically outperforms related neural network kernels on a series of benchmarks.

## 1 Kernel methods, deep learning and implicit deep learning

Kernel methods are a classical paradigm for analysing representational capacity, bias, generalisation performance and practical algorithms for nonparametric prediction (Schölkopf et al., 2002). Many classical nonparametric models can be seen as extensions of parametric models (Saunders, 1998; Rasmussen & Williams, 2006, § 2.2) that allow for increased representational capacity while retaining some statistical model-based properties. Examples of model-based qualities may include the smoothness, stationarity or periodicity of the predictor (Duvenaud, 2014, § 2) or the statistical interpretation of the learning procedure (Sollich, 2002; Rasmussen & Williams, § 3), which may be understood by examining the kernel or the loss function (Banerjee et al., 2005, Theorem 4).

Despite early successes of kernel methods, when data is plentiful and/or modelling is hard, over-parameterised and under-regularised deep learning is now seen as the dominant paradigm for practical nonparametric-style prediction (OpenAI et al., 2019; Adiwardana et al., 2020; Rombach et al., 2022). Unlike parametric and classical nonparametric approaches, the architecture and loss functions of many explicit neural networks are driven purely from the perspective of representational power or predictive performance (either empirical (Vaswani et al., 2017) or mathematical (Raghu et al., 2017)) rather than model-based qualities.

A fruitful direction is to analyse deep learning predictors through the reductionist lens of kernel methods through sufficiently well-behaved neural networks in certain large parameter count regimes (Neal, 1995). However, to the best of our knowledge, no current theory describes architectural properties of neural networks in the kernel regime such as choice of activation function, depth and skip connections, in terms of model-based properties. It is desirable to motivate predictive deep learning architectures from a more fundamental, statistical model-based perspective (Rudin, 2019; Efron, 2020) in a kernel regime.

Implicit neural networks are an emerging approach to model-based deep learning, where the layers are defined to implicitly satisfy the solution to a given problem. For example, deep declarative networks

(DDNs) (Gould et al., 2021) solve optimisation problems, deep equilibrium models (DEQs) (Bai et al., 2019) solve fixed point (algebraic) problems and neural ODEs solve differential equations (Chen et al., 2018). Such problems are usually computed numerically via the (approximate) fixed point of an iterative procedure. This leads to the view that implicit layers are themselves a composition of infinitely many functions. Owing to the complexity of deep learning algorithms, theory falls short of explaining the empirically demonstrated successes of both implicit and explicit models. To the best of our knowledge, no general notion of an implicit kernel is currently described in the literature.

### 1.1 Our contribution: an implicit kernel and an update rule in kernel space

**Updates in feature space** Solutions to optimisation, fixed point or differential equation problems are in practice most often obtained via a possibly stochastic iterative update procedure. Let  $X_1 \in \mathbb{X} \subseteq \mathbb{R}^l$  be an input to the problem and  $\psi_{X_1} \in \boldsymbol{\psi} \subseteq \mathbb{R}^m$  be the solution to the problem. Note that  $\psi_{X_1}$  is the evaluation of an implicit function of  $X_1$ . Let  $\psi_{X_1}^{(t)}$  be a representation of the solution obtained at iterate  $t$ . That is, there exists some possibly stochastic function  $g^{(t)}(\cdot; X) : \boldsymbol{\psi} \rightarrow \boldsymbol{\psi}$  such that

$$\psi_{X_1}^{(t+1)} = g^{(t)}(\psi_{X_1}^{(t)}; X_1), \quad \text{and} \quad \psi_{X_2}^{(t+1)} = g^{(t)}(\psi_{X_2}^{(t)}; X_2). \quad (1)$$

We call  $\psi_{X_1}^{(t)}$  and  $\psi_{X_2}^{(t)}$  features and  $g^{(t)}$  the update rule in feature space. We emphasise that we consider the problem where features are updated, not weight parameters as in some other settings.

**Deep equilibrium kernels** We find helpful the notion of an implicitly defined kernel, which we call a deep equilibrium kernel (DEK). This allows us to draw parallels between infinitely wide implicit neural networks and implicitly defined kernel machines. We consider three kernel evaluations in terms of the implicit updates in feature space,

$$\underbrace{\overline{\Psi}_{12}^{(t+1)} \triangleq \psi_{X_1}^{(t+1)\top} \psi_{X_2}^{(t+1)}}_{\text{finite feature DEK (ffDEK)}}, \quad \underbrace{\Psi_{12}^{(t+1)} \triangleq \text{plim}_{m \rightarrow \infty} \overline{\Psi}_{12}^{(t+1)}}_{\text{DEK}}, \quad \text{and} \quad \underbrace{\Psi_{12} \triangleq \lim_{\tau \rightarrow \infty} \Psi_{12}^{(\tau)}}_{\text{limiting DEK (\ellDEK)}}, \quad (2)$$

where defined, where plim denotes convergence in probability. Note the order of the limits. We will similarly write  $\Psi_{11}$  and  $\Psi_{22}$  to represent evaluations of such DEKs at  $(X_1, X_1)$  and  $(X_2, X_2)$  respectively. We write  $\Psi^{(t+1)}$  for the corresponding  $2 \times 2$  PSD matrices containing  $\Psi_{11}^{(t+1)}$ ,  $\Psi_{12}^{(t+1)}$  and  $\Psi_{22}^{(t+1)}$  (and likewise for the ffDEK and  $\ell$ DEK). The dimensionality of features are allowed to grow to infinity, but only after taking the inner product — resulting in a scalar value for examination — avoiding the necessity of describing, analysing and building algorithms involving infinite dimensional feature spaces.

**Updates in kernel space** Let  $\mathbb{S}_+^2 = \{K \in \mathbb{R}^{2 \times 2} \mid K \succeq 0\}$  denote the space of  $2 \times 2$  PSD matrices. Our central questions are as follows. Firstly, as  $m \rightarrow \infty$ , does there exist a  $G$  such that updates may be performed on  $2 \times 2$  PSD DEK matrices instead of in  $2m$ -dimensional feature space? Secondly, can we write a closed form for  $G$ ? Finally, does this iteration converge? That is, does there exist a closed-form update rule in kernel space  $G(\cdot; X_1, X_2) : \mathbb{S}_+^2 \rightarrow \mathbb{S}_+^2$  such that

$$\Psi^{(t+1)} = G(\Psi^{(t)}; X_1, X_2)? \quad (3)$$

$$\text{And does} \quad \Psi = G(\Psi) = \lim_{\tau \rightarrow \infty} \underbrace{G \circ \dots \circ G}_{\tau \text{ compositions}}(\Psi^{(0)}; X_1, X_2)? \quad (4)$$

For convenience, we notationally decompose  $G$  into components via a function  $G$  satisfying for each  $ij \in \{11, 22, 12\}$

$$G(\Phi; X_1, X_2) = \begin{pmatrix} G_{11} & G_{12} \\ G_{12} & G_{22} \end{pmatrix}, \quad \text{where} \quad (5)$$

$$G_{ij} \triangleq G(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}; X_i, X_j) \triangleq \left( G(\Phi; X_1, X_2) \right)_{ij}.$$

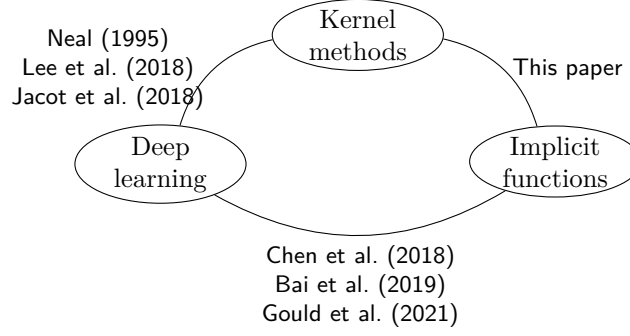


Figure 1: We establish links between kernel methods and implicit functions to design a neural network kernel with corresponding statistical assumptions.

**Contributions** We study an important special case of a DEK where we answer (3) and (4) positively, one in which the features are iteratively updated using SGD applied to a latent variable model. The model is *over-parameterised* (the number of parameters grows much faster than the amount of data), but shallow (the *motivation* for the model more closely resembles exponential family PCA than a deep neural network). The objective (17) to which we apply SGD is an *under-regularised* variant of an expected negative log posterior. Our main result (stated precisely in § 3) is a constructive proof of the existence of  $\mathbf{G}$ .

Surprisingly, despite the feature space of the DEK having seemingly no direct relation with deep learning predictors, deep learning structures emerge as part of our analysis. Our DEK may be understood as an infinitely wide DEQ whose iterates are computed with stochastic approximation rather than as a deterministic fixed point iteration. Further, the kernel iterates of our DEK resemble previously introduced NNKs and NTKs and may be computed with deterministic fixed point solvers.

**Theorem 4 and Corollary 5 (informal).** *When the features  $\psi_{X_1}^{(t)}$  and  $\psi_{X_2}^{(t)}$  are point estimates obtained by SGD applied to objective (17), we construct a deterministic update rule for the DEK,  $\mathbf{G}$ . Repeated applications of  $\mathbf{G}$  converge to a fixed point, the  $\ell$ DEK. In other words, for a specific continuous latent variable model, under mild assumptions, (3) and (4) hold for the kernel update rule  $\mathbf{G}$  on  $2 \times 2$  PSD matrices.*

We further quantify the degree to which the  $\ell$ DEK is an invariant of SGD when treated as an  $\mathbf{f}$ DEK (Theorem 8).

## 2 Background

Our analysis requires combining fixed point theory (optimisation), the exponential family (statistics), and neural network kernels (machine learning). We briefly describe elements of these topics here.

### 2.1 Fixed points and infinite compositions

Let  $f : \mathbb{F} \rightarrow \mathbb{F}$  for some set  $\mathbb{F}$  equipped with a norm  $\|\cdot\|$  and norm-induced metric. A fixed point of  $f$  is any  $Z^* \in \mathbb{F}$  satisfying  $f(Z^*) = Z^*$ . Banach’s fixed point theorem (BFPT) gives sufficient conditions for the existence and uniqueness of such a fixed point.

**Theorem 1 (BFPT).** *Let  $(\mathbb{F}, \|\cdot\|)$  be a non-empty complete normed space. A mapping  $f : \mathbb{F} \rightarrow \mathbb{F}$  is called a contraction mapping if there exists some  $q \in [0, 1)$  such that  $\|f(Z) - f(Z')\| \leq q\|Z - Z'\|$  for every  $Z, Z' \in \mathbb{F}$ . Every contraction mapping  $f$  admits a unique fixed point  $Z^* \in \mathbb{F}$ . Furthermore, for any initial element  $Z^{(1)} \in \mathbb{F}$ , the sequence  $Z^{(t+1)} = f(Z^{(t)})$  for  $t \geq 1$  converges to  $Z^*$  as  $t \rightarrow \infty$ .*

It is worth noting that BFPT not only provides a mathematical condition for well-posedness, but also describes an algorithm for approximating fixed points of contraction mappings. We call this algorithm

the *naive fixed point solver*, which simply involves applying a  $\tau$ -fold composition of  $f$  to some starting value  $Z^{(1)}$ , with a linear rate of convergence immediate from the definition of contraction mapping, i.e.  $\|Z^* - Z^{(t+1)}\| \leq \frac{q^t}{1-q} \|Z^{(2)} - Z^{(1)}\|$ . Other solvers for fixed point problems are available, many of which are approximate Newton methods for root finding (Kelley, 1995).

Deep equilibrium models (DEQs) (Bai et al., 2019) are neural network predictors constructed of parameterised layers that output the solution to fixed point equations  $f_U(Z^*) = Z^*$ . These layers draw upon earlier works on recurrent backpropagation (Pineda, 1987; Almeida, 1990), leveraging the modern machinery of deep learning architectures, optimisers and heuristics. The unsupervised learning problem for a DEQ is

$$\underbrace{\min_U \sum_{i=1}^N L(X_i, Z_i^*, U)}_{\text{Empirical risk minimisation for parameters } U} \quad \text{subject to} \quad \underbrace{Z_i^* = f_U(Z_i^*, X_i)}_{\text{Fixed point solution for DEQ predictions } Z_i^*},$$

where  $L$  is some loss function,  $U$  is a parameter object, and  $\{X_i\}_{i=1}^N$  is a collection of input examples. (A supervised setting might also involve a set of output examples). Derivatives  $\frac{\partial Z_i^*}{\partial U}$  of outputs of these layers with respect to their parameters  $U$  can be computed without backpropagating through the iterates of the fixed point solver using the implicit function theorem (Bai et al., 2019). This allows first-order stochastic gradient methods that are popular with explicit deep learning architectures to be applied to DEQs.

In general, it is not guaranteed that a function necessarily admits a unique fixed point; various works discuss dealing with multiple fixed points or ensuring or encouraging that exactly or at least one fixed point exists (Winston & Kolter, 2020; Revay et al., 2020; El Ghaoui et al., 2021). Interestingly, if a single DEQ layer involves finding the fixed point of a contraction mapping, by Theorem 1 the output computed by a DEQ has the interpretation of an infinitely deep neural network with shared parameters in each layer. More generally, mappings computed by the naive fixed point solver have interpretations as very deep neural networks with shared parameters in each layer. Since zeros of the gradient of sufficiently well behaved objectives are stationary points of the objectives, DEQ layers share a connection with optimisation-based implicit layers (Gould et al., 2021), as explored in various works (Revay et al., 2020; Xie et al., 2021; Tsuchida et al., 2022; Riccio et al., 2022; Tsuchida & Ong, 2022). Our current investigation concerns a connection more specific than optimisation, since it considers the special case of applying SGD.

## 2.2 Exponential families

**Exponential families** The feature mappings that we use to build our kernel are estimates obtained using SGD applied to certain variants of exponential family likelihoods and Gaussian priors. We now define *minimal and regular exponential families in canonical form*. Let  $h$  be a probability density (mass) function supported on data space  $\mathbb{Y} \subseteq \mathbb{R}$ . Let  $T : \mathbb{Y} \rightarrow \mathbb{R}$  be a function called the *sufficient statistic*. Given some *canonical parameter*  $\eta$  belonging to an open set  $\mathbb{H} \subseteq \mathbb{R}$ , we may construct a probability density (mass) function by normalising the nonnegative function  $h(\cdot) \exp(T(\cdot)\eta)$ . The normalising constant is called the partition function, and its strictly convex and infinitely differentiable logarithm  $A$  is called the *log partition function* (Wainwright et al., 2008, Proposition 3.1). We write

$$p(y | \eta) = h(y) \exp(T(y)\eta - A(\eta)), \quad A(\eta) = \log \int_{\mathbb{Y}} h(y) \exp(T(y)\eta) dy$$

for the evaluation of a probability density (mass) function of an exponential family. The log partition function  $A$  acts as a cumulant generating function for the conditional distribution of the sufficient statistic  $T$ . In particular the expected value of the sufficient statistic (often called the *expectation parameter* (Nielsen & Garcia, 2009)) is the gradient of the log partition function  $A$ . That is,

$$\mathbb{E}[T(y) | \eta] = A'(\eta). \quad (6)$$

We consider factorised exponential families in the following sense. Let  $y_1, \dots, y_d$  be distributed according to the same exponential family and define data vector  $Y = (y_1, \dots, y_d)^\top$  and canonical parameter vector  $H = (\eta_1, \dots, \eta_d)$ . Then the joint distribution of data  $Y$  conditioned on canonical parameters  $H$  is the product of the individual elements

$$p(Y | H) = \prod_{i=1}^d p(y_i | \eta_i) = \left( \prod_{r=1}^d h(y_r) \right) \exp \left( T(Y)^\top H - A(H)^\top \mathbf{1} \right), \quad (7)$$

where we write  $T(Y) = (T(y_1), \dots, T(y_d))^\top$ ,  $A(H) = (A(\eta_1), \dots, A(\eta_d))^\top$  and  $\mathbf{1} = (1, \dots, 1)^\top$ .

**Link functions and canonical link functions** Exponential families are used in generalised linear models (GLMs) (McCullagh & Nelder, 1989). In GLMs, the conditional expectation (6) of an exponential family is set to be the result of applying an (invertible) *inverse link function*  $s^{-1}$  to the result of a linear transformation of features  $\phi \in \mathbb{R}^m$  (classically called parameters). That is, for some linear basis  $V \in \mathbb{R}^{d \times m}$  (classically called covariates),

$$A'(H) = \mathbb{E}[T(Y) | H] = s^{-1}(V\phi). \quad (8)$$

The conditional expectation is then mapped to the canonical parameter  $H$  through  $H = (A')^{-1} \circ s^{-1}(V\phi)$ , noting that  $A'$  is invertible because  $A$  is strictly convex. In the case where  $s^{-1}$  is chosen to be  $A'$ ,  $s \equiv (A')^{-1}$  is called the *canonical link function*, and we observe from (8) that the canonical parameter and conditional expectation satisfy

$$H = V\phi, \quad \mathbb{E}[T(Y) | V\phi] = A'(V\phi) = s^{-1}(V\phi). \quad (9)$$

There are two main and sometimes conflicting reasons why one might be interested in using a non-canonical link function. The first is computational; if the link function were canonical, for some distributions such as Gamma or exponential one would need a constrained optimisation method over the open set  $\mathbb{H}$  instead of  $\mathbb{R}$ . If  $s^{-1}$  were allowed to be non-canonical — that is, we are free to choose  $s^{-1}$  different from  $A'$  — we could map the conditional expectation to the appropriate constraint set and unconstrained optimisation procedures could be applied. In a Bayesian context, sampling from the posterior over  $\phi$  can be made easier by convenient choices of  $s$ . For example, the probit model admits an efficient Gibbs sampler for the posterior (Albert & Chib, 1993). The second, and arguably more important consideration is modelling; we might have reason to suspect that the conditional expectation is constrained. For example, if the observations should have a positive expectation, the power family of link functions might be used (McCullagh & Nelder, 1989, equation 2.9a). Alternative link functions can lead to exploiting particular properties of interest; for example, Wiemann et al. (2021) use the softplus function for positive conditional expectations to exploit its identity-like behaviour at large positive values. In weighing up the possibly conflicting aims of computational convenience and modelling suitability, we highlight the view of Efron & Hastie (2021, page 68); while classical exponential families and link functions may lead to closed-form expressions, modern computer technology allows us more flexible models.

**Point estimation** When using a canonical inverse link function  $s \equiv A'$ , the negative logarithm of the likelihood (7) is strictly convex in  $H$ , since  $A$  is strictly convex and linear functions are convex. If  $H$  is chosen to be  $H = V\phi$ , this translates to convexity in  $\phi$ , and maximum likelihood estimates can be computed using first or (more typically, in a classical setting) second order optimisation methods. When  $s^{-1}$  is not a canonical inverse link function, convexity does not necessarily hold. Nevertheless, local estimates are practically useful, so pre-implemented link functions and the option to implement custom link functions is available in a number of software frameworks including **R** (R Core Team, 2021, `family`) and **Stata** (Hardin & Hilbe, 2018, `glm`).

### 2.3 Kernels arising from neural networks

Our main result describes the DEK update rule as a composite function involving evaluations of kernels of a particular form. These are kernels that are constructed from neural network models. In this section, we describe such kernels.

The neural network kernel was first investigated as the covariance function of a certain neural network with random parameters and a single hidden fully connected layer (Neal, 1995). Under mild conditions, as the width of the hidden layer goes to infinity the neural network converges to a Gaussian process. This analysis has since been extended to handle multiple layers (Matthews et al., 2018; Lee et al., 2018), other layer types including convolutional layers (Mairal et al., 2014; Garriga-Alonso et al., 2018; Novak et al., 2019; Yang, 2019a;b), and training under gradient flow via the neural tangent kernel (NTK) (Jacot et al., 2018). Since our motivation is better described in terms of inner products of the features, we favour the view of the neural network kernel as an inner product in an infinitely wide hidden layer rather than a covariance function of a Gaussian process. We note that connections between Bayesian Gaussian processes and kernel methods exist (Kanagawa et al., 2018) and apply to some but not all infinitely wide neural networks.

**Neural network kernel, single hidden layer** Let  $\mathbf{W}^{(1)} \in \mathbb{R}^{d \times n}$  be the weights of a fully connected hidden layer with activation function  $\zeta$  defined over the reals. Suppose each entry of  $\mathbf{W}^{(1)}$  is i.i.d. with distribution  $\mathcal{N}(0, 1)$ <sup>1</sup>. Given an input feature  $\phi_1 \in \mathbb{R}^{n \times 1}$  (we take the convention that vectors are column vectors), the signal in the hidden layer is  $h^{(1)} \triangleq \zeta(\mathbf{W}^{(1)}\phi_1)$ <sup>2</sup>. Here and throughout the paper the symbol  $\triangleq$  means that the object on the left hand side is defined to be the expression on the right hand side. By a strong law of large numbers, a suitably normalised inner product in the hidden layer converges almost surely as  $d \rightarrow \infty$  to an expectation,

$$\frac{1}{d} h_1^{(1)\top} h_2^{(1)} = \frac{1}{d} \zeta(\mathbf{W}^{(1)}\phi_1)^\top \zeta(\mathbf{W}^{(1)}\phi_2) \xrightarrow{a.s.} \mathbb{E}_W [\zeta(W^\top \phi_1) \zeta(W^\top \phi_2)],$$

assuming the right hand side is finite, since the inner product is a sum of i.i.d. random variables. Here  $W^\top \in \mathbb{R}^{1 \times m}$  is a vector with i.i.d. entries drawn from  $\mathcal{N}(0, 1)$ . We define

$$k_\zeta(\phi_1, \phi_2) \triangleq \mathbb{E}_W [\zeta(W^\top \phi_1) \zeta(W^\top \phi_2)], \quad (10)$$

and call  $k_\zeta$  a single hidden layer neural network kernel (NNK) with activation function  $\zeta$ . The PSD kernel  $k_\zeta$  uniquely defines an RKHS by the Moore–Aronszajn theorem. Closed-form expressions of  $k_\zeta$  for different  $\zeta$  are available (Williams, 1997; Le Roux & Bengio, 2007; Cho & Saul, 2009; Tsuchida et al., 2018; Pearce et al., 2019; Tsuchida, 2020; Meronen et al., 2020; Tsuchida et al., 2021; Han et al., 2022).

Define  $(\chi_1, \chi_2)^\top \triangleq (W^\top \phi_1, W^\top \phi_2)^\top$ , which is a zero mean bivariate Gaussian with a covariance matrix  $\Sigma^{(1)}$ . Note that  $k_\zeta(\phi_1, \phi_2)$  depends on the input features  $\phi_1$  and  $\phi_2$  only through the covariance matrix  $\Sigma^{(1)}$ . It is helpful to explicate this dependence structure through a special notation. We have that (10) is equal to

$$\kappa_\zeta(\Sigma_{11}^{(1)}, \Sigma_{22}^{(1)}, \Sigma_{12}^{(1)}) \triangleq k_\zeta(\phi_1, \phi_2) = \mathbb{E}_{(\chi_1, \chi_2)^\top \sim \mathcal{N}(\mathbf{0}, \Sigma^{(1)})} [\zeta(\chi_1) \zeta(\chi_2)], \quad \Sigma^{(1)} \triangleq \begin{pmatrix} \phi_1^\top \phi_1 & \phi_1^\top \phi_2 \\ \phi_2^\top \phi_1 & \phi_2^\top \phi_2 \end{pmatrix}. \quad (11)$$

With an abuse of terminology, we refer to both  $k_\zeta$  and  $\kappa_\zeta$  as PSD single hidden layer NNKs. For a more detailed description of the NNK, see Appendix C.

**Neural network kernel,  $\tau$  hidden layers** One may compose (11) multiple times by applying a sequence of kernels to a 3-dimensional state represented by a  $2 \times 2$  PSD matrix  $\Sigma^{(t)}$ , in place of the infinitely wide signals. This 3-dimensional state represents the two squared norms and inner product in each hidden layer. For  $t = 1, \dots, \tau$  and  $ij \in \{11, 22, 12\}$ ,

$$\Sigma_{ij}^{(t+1)} \triangleq \mathbb{E}_{(\chi_i, \chi_j)^\top \sim \mathcal{N}(\mathbf{0}, \Sigma^{(t)})} [\zeta(\chi_i) \zeta(\chi_j)] = \kappa_\zeta(\Sigma_{ii}^{(t)}, \Sigma_{jj}^{(t)}, \Sigma_{ij}^{(t)}), \quad (12)$$

where  $\Sigma_{ij}^{(t)}$  denotes the  $ij$ th element of  $\Sigma^{(t)}$ . This iteration appears in deep infinitely wide NNKs (Matthews et al., 2018; Lee et al., 2018). We will refer to this kernel as the  $\tau$  layer NNK. Evaluations  $\Sigma_{12}^{(\tau+1)}$  of the PSD kernel are determined entirely by the activation function  $\zeta$ , and uniquely define an RKHS.

<sup>1</sup>The effect of non-unit weight variance may be obtained by scaling all inputs  $\phi_1$  by a hyperparameter. Similarly, arbitrary covariance structures inside rows of  $\mathbf{W}^{(1)}$  can be reflected as linear transformations of all inputs  $\phi_1$ .

<sup>2</sup>The effect of zero mean Gaussian biases may be obtained by augmenting inputs with an additional coordinate. The magnitude of this coordinate is equivalent to the quotient of the standard deviation of the weights to the biases.

**Neural tangent kernel,  $\tau$  hidden layers** This kernel describes the limiting behaviour of randomly initialised neural networks that are trained under gradient flow (Jacot et al., 2018). The kernel iterations are similar to (12), but also contain components involving the derivative  $\dot{\zeta}$  of  $\zeta$ . Let  $\odot$  denote elementwise product. In addition to the iteration (12), let  $\Theta^{(1)} = \Sigma^{(1)}$  and define

$$\Theta^{(t+1)} \triangleq \Theta^{(t)} \odot \dot{\Sigma}^{(t+1)} + \Sigma^{(t+1)}, \quad \text{where} \quad (13)$$

$$\dot{\Sigma}_{ij}^{(t+1)} \triangleq \mathbb{E}_{(\chi_i, \chi_j)^\top \sim \mathcal{N}(\mathbf{0}, \Sigma^{(t)})} [\dot{\zeta}(\chi_i) \dot{\zeta}(\chi_j)] = \kappa_{\dot{\zeta}}(\Sigma_{ii}^{(t)}, \Sigma_{jj}^{(t)}, \Sigma_{ij}^{(t)})$$

to obtain the evaluation of the PSD NTK in the last iteration  $\Theta_{12}^{(\tau+1)}$ . Once again, the kernel is determined entirely by the  $\zeta$ , and uniquely defines an RKHS.

## 2.4 Notation

Numerical subscripts are used to extract (groups of) indices of a vector or matrix. Parenthesised superscripts indicate a layer or iteration of a naive fixed point solver, both of which turn out to be the same in our constructions. We index objects by iteration by superscript  $(t)$ , so that  $\psi^{(t)}$  represents a feature in the  $t$ th iteration. We use  $\phi$  and  $\Phi$  for arbitrary vectors and inner products that are not necessarily obtained by iterations of SGD. We will use  $\psi$  and  $\Psi$  for feature mappings and inner products of feature mappings that are obtained by iterations of SGD.

We assume that we are given access to a dataset  $\mathbf{X} \in \mathbb{R}^{N \times l}$  of  $N$  examples of datapoints  $X_i \in \mathbb{X} \subseteq \mathbb{R}^l$ . We denote by  $X_1$  and  $X_2$  any two elements of this dataset.

There are two types of function signatures we associate with PSD kernels. The first is for a usual PSD kernel  $k : \mathbb{X} \times \mathbb{X} \rightarrow \mathbb{R}$ , so that an evaluation is written  $k(X, X')$  for any two  $X, X' \in \mathbb{X}$ . We call this form a  $k$ -form kernel. The second is for a PSD kernel whose evaluation depends on  $\phi_1, \phi_2 \in \boldsymbol{\psi}$  only through evaluations  $\Phi_{11} = \langle \phi_1, \phi_1 \rangle, \Phi_{12} = \langle \phi_1, \phi_2 \rangle, \Phi_{22} = \langle \phi_2, \phi_2 \rangle$  of some suitably defined inner product  $\langle \cdot, \cdot \rangle : \boldsymbol{\psi} \times \boldsymbol{\psi} \rightarrow \boldsymbol{\psi}$ . We represent such a kernel through  $\kappa : \boldsymbol{\psi}^3 \rightarrow \mathbb{R}$  with evaluations  $\kappa(\langle \phi_1, \phi_1 \rangle, \langle \phi_2, \phi_2 \rangle, \langle \phi_1, \phi_2 \rangle)$ . An example of this second form is the NNK (11). We call this form a  $\kappa$ -form kernel.

Our notation is summarised in Table 3 in Appendix A.

## 3 Main results

Our results are most clearly described in terms of an alternative parameterisation of exponential families and link functions, which are perhaps of independent interest. We first describe this alternative parameterisation in § 3.1, before moving onto the setup for our main analysis in § 3.2. We then in § 3.3 provide a special case (Corollary 3) of our main and most general result (Theorem 4) in § 3.4. Finally, quantification of error between the DEK and ffDEK is described in § 3.5. We give examples of our resulting updates in Appendix G.2.

### 3.1 An alternative view of link functions in exponential families

Instead of computing via the conditional expectation resulting from the application of an inverse link function (8), we follow Tsuchida & Ong (2022) and learn the canonical parameter via a nonlinearity  $H = R(\mathbf{V}\phi)$ , for some once-differentiable  $R : \mathbb{R} \rightarrow \mathbb{H}$  called the *canonical nonlinearity*. This means that the conditional likelihood (7) is now

$$p(Y | \mathbf{V}, \phi) = \prod_{i=1}^d p(y_i | \eta_i) = \left( \prod_{r=1}^d h(y_r) \right) \exp \left( T(Y)^\top R(\mathbf{V}\phi) - A(R(\mathbf{V}\phi))^\top \mathbf{1} \right). \quad (14)$$

Such a parameterisation is rich enough to recover the (non-canonical) inverse link function view of the statistician (see Proposition 2). It can therefore be considered to be a change of notation, placing emphasis on the canonical nonlinearity  $R$  instead of the inverse link function  $s^{-1}$ . In our setting, one

Exponential family	$A(\eta)$	$s^{-1}(a)$	$R(a)$	$\rho(a)$	$\sigma(a)$
Gaussian	$\eta^2/2$	$s^{-1}(a)$	$s^{-1}(a)$	$(s^{-1})'(a)$	$s^{-1}(a)(s^{-1})'(a)$
Gaussian	$\eta^2/2$	$a$	$a$	1	$a$
Gaussian	$\eta^2/2$	$\text{erf}(a/\sqrt{2})$	$\text{erf}(a/\sqrt{2})$	$2p(a)$	$\text{erf}(a/\sqrt{2})2p(a)$
Gaussian	$\eta^2/2$	$\text{ReLU}(a)$	$\text{ReLU}(a)$	$u(a)$	$\text{ReLU}(a)$
Poisson	$\exp(\eta)$	$s^{-1}(a)$	$\log s^{-1}(a)$	$\frac{(s^{-1})'(a)}{s^{-1}(a)}$	$(s^{-1})'(a)$
Poisson	$\exp(\eta)$	$\exp(a)$	$a$	1	$\exp(a)$
Poisson	$\exp(\eta)$	$\log(1 + \exp a)$	$\log \log(1 + \exp a)$	$\frac{\exp(a)}{(1 + \exp a) \log(1 + \exp a)}$	$\exp(a)/(1 + \exp(a))$
Bernoulli	$\log(1 + \exp(\eta))$	$s^{-1}(a)$		$\frac{(s^{-1})'(a)}{s^{-1}(a)(1 - s^{-1}(a))}$	$\frac{(s^{-1})'(a)}{1 - s^{-1}(a)}$
Bernoulli	$\log(1 + \exp(\eta))$	$\exp(a)/(1 + \exp(a))$	$a$	1	$\exp(a)/(1 + \exp(a))$
Bernoulli	$\log(1 + \exp(\eta))$	$P(a)$	$\log\left(\frac{P(a)}{1 - P(a)}\right)$	$\frac{p(a)}{P(a)P(-a)}$	$\frac{p(a)}{P(-a)}$

Table 1: These examples are obtained by plugging the desired log partition function  $A$  and inverse link function  $s^{-1}$  into expressions (23), (24) and (25). Canonical link functions are shown in blue. General inverse link function settings are shown in red. Here  $P$  and  $p$  respectively denote the cdf and pdf of the univariate standard Gaussian and erf denotes the error function.

advantage of such a notation is that it avoids more complicated function compositions involving inverses and derivatives. For example, instead of writing  $A \circ (A')^{-1} \circ s^{-1}(\mathbf{V}\phi)$  we may write  $A \circ R(\mathbf{V}\phi)$ . The value of these simple compositions become more evident in Proposition 2.

**Nonlinearities and activation functions** The derivatives of the log likelihood (14) (a *score function*) play a central role in numerical procedures associated with estimation. In our setting, such derivatives involve terms derived from  $A$  and  $R$ . These terms are expressed in terms of functions we call *factor activations*  $\rho(a) \triangleq R'(a)$  and *chain activations*  $\sigma(a) \triangleq (A \circ R)'(a)$ . The following identities show how one may map between choices of  $(A, s)$  and choices of  $(A, R)$ , and additionally how these induce activation functions  $\sigma$  and  $\rho$  which appear in gradient-based optimisers and our later derivations. Note that we may choose the inverse link function  $s^{-1}$  to be non-canonical (not  $A'$ ).

**Proposition 2.** *Consider a regular and minimal exponential family with log partition function  $A : \mathbb{H} \rightarrow \mathbb{R}$ . Suppose the conditional expectation belongs to a set  $\mathbb{A}'$ , that is,  $A'(\eta) \in \mathbb{A}'$  for all canonical parameters  $\eta \in \mathbb{H}$ . Let  $s^{-1} : \mathbb{B} \rightarrow \mathbb{A}'$  be an inverse link function, for some  $\mathbb{B} \subseteq \mathbb{R}$ . That is, for every  $\eta \in \mathbb{H}$  there exists some  $a \in \mathbb{B}$  such that  $A'(\eta) = s^{-1}(a)$ . Then equivalently,  $\eta = R(a)$ , where  $R : \mathbb{B} \rightarrow \mathbb{H}$  is defined by  $R(a) \triangleq ((A')^{-1} \circ s^{-1})(a)$ . Furthermore,*

$$\rho(a) \triangleq R'(a) = \frac{(s^{-1})'(a)}{A'' \circ (A')^{-1} \circ s^{-1}(a)} \quad \text{and} \quad \sigma(a) \triangleq (A \circ R)'(a) = \frac{s^{-1}(a)(s^{-1})'(a)}{A'' \circ (A')^{-1} \circ s^{-1}(a)}.$$

The proof is given in Appendix B. In practice we will take  $\mathbb{B} = \mathbb{R}$ . We observe that  $R$  is the identity if and only if  $s$  is a canonical link function (which is to say that  $s^{-1}(a) = A'(a)$ ). For the special case of a Gaussian with known variance,  $A'$  is the identity and  $R$  is the inverse link function  $s^{-1}$ . Further cases are listed in Table 1.

Nonlinear parameterisation framed in terms of  $R$  instead of  $s^{-1}$  are often used (McCullagh & Nelder, 1989, Chapter 11.4 and references therein), but their general relationship to  $s^{-1}$  does not appear to be discussed. As our setting is equivalent to using an arbitrary link function, we inherit the motivation of using a non-identity  $R$  from the motivation for using a non-canonical link function. We also inherit the usual difficulties in estimation and sample complexity due to using not necessarily canonical link functions.

Recall that the choice of  $(A, s)$  should be informed by both modelling and numerical convenience (sampling, optimisation) considerations. Motivated by neural network kernels, we find a different set of  $(A, s)$  pairs convenient to work with compared with the generalised linear model setting. Convenience here translates to being able to compute certain Gaussian integrals of the form (11) in closed form. For example, we find it easy to work with a Gaussian with non-negative conditional expectation, parameterised by  $A(\eta) = \eta^2/2$  and  $s^{-1}(a) = \text{ReLU}(a)$ , where ReLU is the popular rectified linear unit (Efron



& Hastie, 2021, page 362). Another convenient setting is a Gaussian likelihood and probit inverse link function (in contrast with the often seen Bernoulli and probit inverse link function). Our theory holds for general  $(A, s)$  pairs, but its practical efficiency is contingent upon the existence of efficient numerical routines for computing the integral (11). Such numerical routines in the absence of closed-forms are available in other works (Zandieh et al., 2021; Han et al., 2022), but we do not study their application here.

### 3.2 Setup

**Stochastic gradient descent** We apply SGD (Wright & Recht, 2022, Chapter 5) to a minimisation objective  $\mathcal{L}(\phi; X) = \mathbb{E}_{\mathbf{V}} L(\phi; X, \mathbf{V})$  with decision variable  $\phi$ , input  $X$  and random object  $\mathbf{V}$ . Given two inputs  $X_1$  and  $X_2$ , the  $t + 1$ th iterates are

$$\psi_{X_1}^{(t+1)} = \psi_{X_1}^{(t)} - \alpha^{(t)} \frac{\partial}{\partial \psi_{X_1}^{(t)}} L(\psi_{X_1}^{(t)}; X_1, \mathbf{V}^{(t)}) \quad \text{and} \quad \psi_{X_2}^{(t+1)} = \psi_{X_2}^{(t)} - \alpha^{(t)} \frac{\partial}{\partial \psi_{X_2}^{(t)}} L(\psi_{X_2}^{(t)}; X_2, \mathbf{V}^{(t)}) \quad (15)$$

with initial features  $\psi_{X_1}^{(0)}$  and  $\psi_{X_2}^{(0)}$ , a sequence  $\{\alpha^{(t)}\}_t$  of step sizes, and a sequence of  $\{\mathbf{V}^{(t)}\}_t$  of iid samples of  $\mathbf{V}$ . We use the features  $\psi_{X_1}^{(t)}$  and  $\psi_{X_2}^{(t)}$  in (15) to define the ffDEK, DEK and  $\ell$ DEK via (2). We stress that we are updating features, not weight parameters.

**Continuous latent variable model** We work with a data generating process which is a slight nonlinear generalisation (Tsuchida & Ong, 2022) of exponential family PCA (Collins et al., 2001), allowing for nonlinear  $R$  (or equivalently, non-canonical link functions as in Proposition 2). This model describes data  $Y$  as being drawn from an exponential family distribution with a canonical parameter that is a function of a latent  $\phi$ .

More concretely, let  $X \in \mathbb{X} \subseteq \mathbb{R}^l$  be an input and suppose that data  $Y = \Gamma(X)$  follows a factorised exponential family (7) for some realisation of a random mapping  $\Gamma : \mathbb{X} \rightarrow \mathbb{Y}^d$ . Let  $R : \mathbb{R} \rightarrow \mathbb{H}$  be a once-differentiable function. Choose the canonical parameter  $H = R(\mathbf{V}\phi)$  to be the composition of  $R$  and a linear transformation  $\mathbf{V}$  of a latent input variable  $\phi \in \Psi = \mathbb{R}^m$ . Place an i.i.d.  $\mathcal{N}(0, 1)$  prior over each entry of  $\mathbf{V} \in \mathbb{R}^{d \times m}$ , independent of  $\Gamma$ . Place an i.i.d.  $\mathcal{N}(0, \lambda^{-1}/m)$  prior over  $\phi$ . This results in a pre-nonlinearity parameter  $\mathbf{V}\phi$  having components with variance which stays in  $d$  and  $m$ . For some constant  $C$  not depending on  $\phi$ , we have

$$-\log p(\phi \mid \Gamma(X), \mathbf{V}) = - \left( \underbrace{\log p(\Gamma(X) \mid R(\mathbf{V}\phi))}_{\text{Log likelihood}} - \underbrace{m \frac{\lambda}{2} \|\phi\|^2}_{\text{-Log prior}} \right) + C.$$

As discussed in Proposition 2, the derivative of the negative log-posterior with log-partition function  $A$ , which appears in our later optimisation procedure, induces two functions  $\rho(a) \triangleq R'(a)$  and  $\sigma(a) \triangleq (A \circ R)'(a)$  which we call factor activations and chain activations respectively.

**Objective function** The expected negative log posterior

$$\bar{\mathcal{L}}(\phi; X) \triangleq \mathbb{E}_{\mathbf{V}} \bar{L}(\phi; X, \mathbf{V}), \quad \text{where} \quad \bar{L}(\phi; X, \mathbf{V}) \triangleq \frac{1}{d} \left( -\log p(\Gamma(X) \mid R(\mathbf{V}\phi)) + \textcolor{red}{m} \frac{\lambda}{2} \|\phi\|^2 \right), \quad (16)$$

is a commonly used minimisation objective to find point estimates of  $\phi$ . See Appendix H.1 for a discussion on this objective. The division by  $d$  is introduced to account for the natural numerical scaling of the likelihood term, which is a sum of  $d$  parts. Following recent deep learning trends, we consider an *over-parameterised* and *under-regularised* variant

$$\mathcal{L}(\phi; X) \triangleq \mathbb{E}_{\mathbf{V}} L(\phi; X, \mathbf{V}), \quad \text{where} \quad L(\phi; X, \mathbf{V}) \triangleq \frac{1}{d} \left( -\log p(\Gamma(X) \mid R(\mathbf{V}\phi)) + \textcolor{blue}{\sqrt{md}} \frac{\lambda}{2} \|\phi\|^2 \right), \quad (17)$$

where  $d < m$ . This expected negative log posterior may be obtained by choosing an overly broad i.i.d. prior  $\mathcal{N}(0, \lambda^{-1}/\sqrt{md})$  over  $\phi$ . We will take  $d$  to be a well-behaved function of  $m$  such that  $d \rightarrow \infty$  as  $m \rightarrow \infty$  (see Assumption 1).

**Assumptions** We now describe two assumptions common to both settings. Our first assumption says that the dimensionality  $m$  of the feature space should become larger much faster than the dimensionality  $d$  of the exponential family.

**Assumption 1.** Consider any limit path  $(m, d) \rightarrow (\infty, \infty)$  such that  $\lim_{m \rightarrow \infty} \frac{d}{m} = 0$ .

Our second assumption describes how the step size  $\alpha$  should depend on  $m$ ,  $d$ , and the SGD iteration  $t$ . Recall the nomenclature “fixed” and “decreasing” as qualifiers for step-size, which describe a dependency on  $t$  (but not  $d$ ). Recall that  $\lambda$  is the regularisation parameter.

**Assumption 2 (a).**  $\lim_{m \rightarrow \infty} \alpha^{(t)} \lambda \sqrt{\frac{m}{d}} = 1$ .

Assumption 2 (a) allows for fixed or decreasing step sizes, such as  $\frac{1}{\lambda} \frac{\sqrt{d}}{\sqrt{m+r(t)}}$  for increasing but finite  $r$ .

We will find that in our setup, the DEK is a composite function involving NNK building blocks.

### 3.3 Error function inverse link and Gaussian likelihood to match a random mapping

We will find that the update rule  $G$  of the DEK is a composite function involving NNK building blocks. In order to clearly highlight the role of these NNK building blocks, we first present a special case of our more general Theorem 4. This provides a clear link between the statistical likelihood model and closed form expressions for the DEK update rule.

We choose an exponential family, Canonical nonlinearity and random mapping  $\Gamma$ . This particular setup leads to closed-form expressions for the NNKs involved in the update rule. As an activation function, we choose the error function  $\text{erf}(z) = 2/\sqrt{\pi} \int_0^z e^{-v^2} dv$  (closely related to the Probit function), and rely on a closed-form NNK derived in Williams (1997),

$$\kappa_{\text{erf}(\cdot/\sqrt{2})}(\Sigma_{11}, \Sigma_{22}, \Sigma_{12}) = \frac{2}{\pi} \sin^{-1} \frac{\Sigma_{12}}{\sqrt{(1 + \Sigma_{11})(1 + \Sigma_{22})}}. \quad (18)$$

We show here the statistical modelling choices and their corresponding effect on the DEK update rule. In this special case, our main result (Theorem 4) implies Corollary 3, as proven in Appendix G.1. Recall from § 2.2, that the designer needs to choose a log partition function  $A$ , and a canonical nonlinearity  $R$ . We map input  $X \in \mathbb{R}^l$  to data  $Y \in \mathbb{R}^d$  through a random mapping  $Y = \text{erf}(WX/\sqrt{2}) + Q$ , where  $W \in \mathbb{R}^{d \times l}$  and  $Q \in \mathbb{R}^d$  contain i.i.d. standard Gaussian elements. The distribution of  $Y$  given  $\text{erf}(WX/\sqrt{2})$  is conditionally Gaussian, with conditional expectation  $\text{erf}(WX/\sqrt{2})$  having elements between  $-1$  and  $1$ . We therefore choose a matching inverse link function, to represent the conditional expectation as a function of features  $\psi_X$ . The inverse link function is  $s^{-1}(a) = \text{erf}(a/\sqrt{2})$ . The log partition function is  $A(\eta) = \eta^2/2$  and the sufficient statistic is  $T(y) = y$ . Since the likelihood is Gaussian, the canonical nonlinearity  $R$  and inverse link function  $s^{-1}$  are the same, as shown in the first row of Table 1. In this particular case, the activations  $\rho$  and  $\sigma$  are shown in the third row of Table 1.

**Corollary 3.** Suppose input  $X$  is mapped to data  $Y$  by  $Y = \text{erf}(WX/\sqrt{2}) + Q$ , where  $\text{erf}$  is the error function and  $W \in \mathbb{R}^{d \times l}$  and  $Q \in \mathbb{R}^d$  contain i.i.d. standard Gaussian elements. Choose the log partition function  $A(\eta) = \eta^2/2$ . Choose the canonical nonlinearity  $R(a) = \text{erf}(a/\sqrt{2})$ , or equivalently, choose the inverse link function to be  $s^{-1}(a) = \text{erf}(a/\sqrt{2})$ . This implies that  $\rho(a) = 2p(a)$  and  $\sigma(a) = 2p(a) \text{erf}(a/\sqrt{2})$ , where  $p$  is the pdf of the standard Gaussian. Then  $\kappa_\rho$  and  $\kappa_\sigma$  are given by

$$\begin{aligned} \kappa_\rho(\Phi_{11}, \Phi_{22}, \Phi_{12}) &= \frac{2}{\pi \sqrt{(1 + \Phi_{11})(1 + \Phi_{22}) - \Phi_{12}^2}}, \\ \kappa_\sigma(\Phi_{11}, \Phi_{22}, \Phi_{12}) &= \kappa_\rho(\Phi_{11}, \Phi_{22}, \Phi_{12}) \kappa_{\text{erf}(\cdot/\sqrt{2})}(F_{11}, F_{22}, F_{12}), \quad \text{where } F = (\Phi^{-1} + I)^{-1}. \end{aligned}$$

Let  $C_{ij} = \kappa_{\text{erf}(\cdot/\sqrt{2})}(X_i^\top X_i, X_j^\top X_j, X_i^\top X_j)$ . Suppose Assumptions 1 and 2 (a) hold. Then applying SGD to objective (17), the update rule  $G$  (3) exists and can be decomposed into  $G$  (5) satisfying

$$G(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}; X_i, X_j) = \frac{1}{\lambda^2} \left( C_{ij} \kappa_\rho(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) + \kappa_\sigma(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) \right).$$

Note that in this case the component  $G$  of the update rule  $\mathbf{G}$  can be computed entirely in closed form. The  $2 \times 2$  matrix  $\mathbf{F} = (\Phi + \mathbf{I})^{-1}$  has a simple closed-form in terms of  $\Phi$  (see Appendix G.1). Recall that the decomposition of  $\mathbf{G}$  into  $G$  says that, by plugging (5) into (3), for each  $ij \in \{11, 22, 12\}$ ,

$$\Psi_{ij}^{(\tau+1)} = G(\Psi_{ii}^{(\tau)}, \Psi_{jj}^{(\tau)}, \Psi_{ij}^{(\tau)}; X_i, X_j). \quad \text{That is, } \Psi^{(\tau+1)} = \mathbf{G}(\Psi^{(\tau)}; X_1, X_2).$$

### 3.4 General case

We now consider the general setting, allowing for arbitrary  $(A, s)$  pairs and random mappings  $\Gamma$ . In order to analyse this generalised setting, we require one additional definition (19) and two additional assumptions 3 and 4.

In the most general setting, the DEK includes some non-symmetric (hence not PSD and not a kernel) cross terms. Given two activations  $\zeta_1$  and  $\zeta_2$ ,

$$\kappa_{\zeta_1, \zeta_2}(\Sigma_{11}, \Sigma_{22}, \Sigma_{12}) \triangleq \mathbb{E}_{(\chi_1, \chi_2)^\top \sim \mathcal{N}(\mathbf{0}, \Sigma)} [\zeta_1(\chi_1) \zeta_2(\chi_2)]. \quad (19)$$

The third assumption says that if the inner products were empirical estimates of an expectation, the resulting expectation is real valued and finite. Recall that  $T$  is the sufficient statistic of the exponential family,  $\Gamma$  is the random mapping from input space to data space, and  $\sigma(a) = (A \circ R)'(a)$

**Assumption 3.** *The expectation  $K(a) = \mathbb{E}_Z \left[ \left( T(\Gamma(X)) \odot \rho(aZ) - \sigma(aZ) \right)^2 \right]$  is finite for all  $X \in \mathbb{X}$  and  $a \in \mathbb{R}$ , where  $Z$  is a standard Gaussian random variable.*

The fourth assumption describes the properties of the random mapping  $\Gamma : \mathbb{X} \rightarrow \mathbb{Y}^d$  as  $d \rightarrow \infty$ . In order to understand what happens to the solutions found by SGD as  $d$  becomes large, we need the inputs which are passed through  $\Gamma$  to be well-behaved. It suffices that a kernel and average defined by  $\Gamma$  converges. We call the limiting kernel  $c$  the *explicit kernel*, which contrasts with our implicitly defined DEK. We give examples in Appendix F.

**Assumption 4.** *The PSD kernel  $c$  defined by  $c(X_1, X_2) \triangleq \lim_{m \rightarrow \infty} \frac{1}{d} T(\Gamma(X_1))^\top T(\Gamma(X_2)) = \mathbb{E} T(\Gamma(X_1))^\top T(\Gamma(X_2))$  is finite. Similarly, the mean function defined by  $\mu(X_1) \triangleq \lim_{m \rightarrow \infty} \frac{1}{d} T(\Gamma(X_1))^\top \mathbf{1} = \mathbb{E} T(\Gamma(X_1))^\top \mathbf{1}$  is finite.*

Our main result is a constructive proof for the existence of an update rule  $\mathbf{G}$ , as posed in (3).

**Theorem 4.** *Suppose Assumptions 1, 2 (a), 3, and 4 hold. Let  $C_{ij} = c(X_i, X_j)$  and  $\mu_i = \mu(X_i)$  be as defined in Assumption 4. Then applying SGD to objective (17), the update rule  $\mathbf{G}$  (3) exists and can be decomposed into  $G$  (5) satisfying*

$$\begin{aligned} & G(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}; X_i, X_j) \\ &= \frac{1}{\lambda^2} \left( C_{ij} \kappa_\rho(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) - \kappa_{\sigma, \rho}(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) \mu_i - \kappa_{\rho, \sigma}(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) \mu_j + \kappa_\sigma(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) \right). \end{aligned}$$

Here  $\kappa_\sigma$ ,  $\kappa_\rho$ ,  $\kappa_{\sigma, \rho}$  and  $\kappa_{\rho, \sigma}$  are as defined by (11), (19) and Proposition 2.

The proof is given in Appendix E. Recall again that the decomposition of  $\mathbf{G}$  into  $G$  says that, by plugging (5) into (3), for each  $ij \in \{11, 22, 12\}$ ,

$$\Psi_{ij}^{(\tau+1)} = G(\Psi_{ii}^{(\tau)}, \Psi_{jj}^{(\tau)}, \Psi_{ij}^{(\tau)}; X_i, X_j). \quad \text{That is, } \Psi^{(\tau+1)} = \mathbf{G}(\Psi^{(\tau)}; X_1, X_2).$$

Note the cross terms involving  $\kappa_{\sigma, \rho}$  and  $\mu$  which were not present in the special case of Corollary 3. These cross-terms arise from random mappings  $\Gamma$  with an average element that is non-zero. In the case of Corollary 3, these cross-terms cancel out.

Theorem 4 implies a fixed point condition by Theorem 1, providing a positive answer for (4).

**Corollary 5.** *Suppose the same setting as Theorem 4. If  $G$  is a contraction mapping, then the DEK converges to a unique fixed point as  $t \rightarrow \infty$ . That is, for each  $ij \in \{11, 22, 12\}$ ,*

$$\Psi_{ij} = \frac{1}{\lambda^2} \left( C_{ij} \kappa_\rho(\Psi_{ii}, \Psi_{jj}, \Psi_{ij}) - \kappa_{\sigma, \rho}(\Psi_{ii}, \Psi_{jj}, \Psi_{ij}) \mu_i - \kappa_{\rho, \sigma}(\Psi_{ii}, \Psi_{jj}, \Psi_{ij}) \mu_j + \kappa_\sigma(\Psi_{ii}, \Psi_{jj}, \Psi_{ij}) \right). \quad (20)$$

Whether  $G$  is a contraction can be determined by a derivative test and an identity given in Theorem 20, as we demonstrate in § G.2. Note that even if a unique fixed point does not exist (which may be the case if  $G$  is not a contraction), one may still compute with finite- $t$  iterates of SGD via Theorem 4.

**Remark 6.** *We may compute iterates of SGD in the limit via the update rule for any  $\tau$  to obtain  $\Psi_{ij}^{(\tau)}$ , which is the naive fixed point algorithm applied to (20). Alternatively, we may compute the  $\ell$ DEK by solving (20) for each  $ij \in \{11, 22, 12\}$  using any other fixed point solver.*

**Notable special cases** Some further examples arising from special choices of  $A$ ,  $R$  and  $\Gamma$  (inducing corresponding  $\rho$ ,  $\sigma$ ,  $c$  and  $\mu$ ) are discussed in Appendix G.2. We find that the linear Gaussian ( $A(\eta) = \eta^2/2$ ,  $R(a) = a$ ) results in a DEK that is a scale multiple of  $c$  (Appendix G.2.1). We can recover an NNK with activation  $\sigma$  when  $A$  and  $R$  are allowed to be general and  $C$  and  $\mu$  are set to zero (Appendix G.2.2). The setting we found useful for our experiments (§ 4) is a nonlinearly parameterised Gaussian ( $A(\eta) = \eta^2/2$  and  $R(a) = \text{ReLU}(a)$ ) with a first-order arc-cosine kernel for  $c$  and a corresponding mean function  $\mu$ . This setting admits a closed-form update rule for  $G$ . See Appendices G.2.4 and G.2.5 for details.

### 3.5 Sensitivity

Assumption 2 (a) suffices for our limiting result to hold. In order to quantify the distance between the limit and finite-dimensional kernels, we require the stronger Assumption 2 (b), which in particular requires a fixed step-size. Both variants 2 (a) and 2 (b) result in a limiting step size of 0, under Assumption 1.

**Assumption 2 (b).** *We have a fixed step-size  $\alpha^{(t)} = \frac{1}{\lambda} \sqrt{\frac{d}{m}}$ .*

Finally, in order to quantify the rate at which the infinite-dimensional, infinite-iteration DEK converges with respect to the dimension, we need the feature mapping to be well-behaved. Lipschitzness and boundedness allows concentration inequalities to be applied.

**Assumption 5.** *Suppose  $\rho$  is bounded or Lipschitz. Suppose  $\sigma$  is bounded or Lipschitz.*

We may quantify the degree to which the infinite-dimensional, infinite-iteration DEK is an invariant of SGD. We define specific values of finite dimensional  $\phi$  and  $\phi'$  using the infinite dimensional kernel fixed point  $\Psi_{11}$ ,  $\Psi_{22}$  and  $\Psi_{12}$ . This definition will serve as a good approximation of an invariant.

**Definition 7.** *Define*

$$r_1 = \sqrt{\Psi_{11}} (1, 0, \dots, 0)^\top \in \mathbb{R}^m, \quad r_2 = \sqrt{\Psi_{22}} (\cos \omega, \sin \omega, 0, \dots, 0)^\top \in \mathbb{R}^m,$$

where  $\cos \omega = \frac{\Psi_{12}}{\sqrt{\Psi_{11}\Psi_{22}}}$ . Then  $r_1^\top r_2 = \Psi_{12}$ ,  $r_1^\top r_1 = \Psi_{11}$ , and  $r_2^\top r_2 = \Psi_{22}$ .

We bound the residual of the finite dimensional kernel evaluated at an initial guess that is the solution of the infinite  $(m, d)$  system. When this bound is small, intuitively speaking, the limiting solution is “almost” an invariant of the finite-dimensional system.

**Theorem 8.** *Suppose Assumptions 1, 2 (b), 3, 4 and 5 hold. Let initial guesses be  $\psi_{X_1}^{(0)} = r_1$  and  $\psi_{X_2}^{(0)} = r_2$  as in Definition 7. Then there exist constants  $Q_2, Q_3, c_2, c_3 > 0$  such that for all  $\delta > 0$ ,  $\epsilon_2 > 0$  and  $\epsilon_2$ ,*

$$\mathbb{P}(|\bar{\Psi}_{12}^{(1)} - \Psi_{12}| \leq \epsilon_1 + \epsilon_2) \geq 1 - \delta_1 - \delta_2,$$

where

$$\varepsilon_1 = \frac{K + \varepsilon_2}{\lambda^2} (2\varepsilon_1 + \varepsilon_1^2), \quad \delta_1 = 2 \exp(-c_2 d M_2) + \exp(-m \delta^2 / 2) \quad \text{and} \quad \delta_2 = 2 \exp(-d c_3 M_3)$$

and  $\varepsilon_1 = \sqrt{\frac{d}{m}} + \delta$ ,  $M_2 = \min\left\{\frac{\varepsilon_2^2}{Q_2^2}, \frac{\varepsilon_2}{Q_2}\right\}$  and  $M_3 = \min\left\{\frac{\varepsilon_2^2}{Q_3^2}, \frac{\varepsilon_2}{Q_3}\right\}$  and  $c_3 > 0$  is some absolute constant.

The proof is given in Appendix E.

## 4 Experiments

Recall that the ffDEK is defined for finite SGD iteration  $\tau$  as an inner product of finite  $m$ -dimensional features. The DEK is defined for finite SGD iteration  $\tau$  as a limit as  $m \rightarrow \infty$  of an inner product of  $m$ -dimensional features. The  $\ell$ DEK is defined as a limit as SGD iteration  $\tau \rightarrow \infty$  of the DEK. Although the DEK and  $\ell$ DEK are defined in terms of infinite dimensional features, evaluations of the DEK and  $\ell$ DEK are scalar values and can be used to form matrices with a finite number of rows and columns. These matrices can be used in downstream kernel algorithms to build predictive algorithms.

### 4.1 Measuring finite-width effects

We empirically measure the similarity of (finite- $\tau$ , finite- $d$ ) ffDEK matrices and (infinite- $\tau$ , infinite- $d$ )  $\ell$ DEK matrices using the centered variant (Cortes et al., 2012) of kernel alignment (Cristianini et al., 2001), abbreviated CKA, as  $d$  increases. We vary  $d$  between 5 and 500 in steps of 5 and choose  $m = d^{3/2}$ . For control, we also measure the CKA between the (finite- $\tau$ , finite- $d$ ) ffDEK and the squared exponential kernel (SEK). See Figure 2, and Appendix J.1 for full details on the experimental setup. As expected (Theorem 4), the CKA between the DEK matrices becomes larger as  $d$  and  $m$  increase, but not between the SEK and finite DEK.

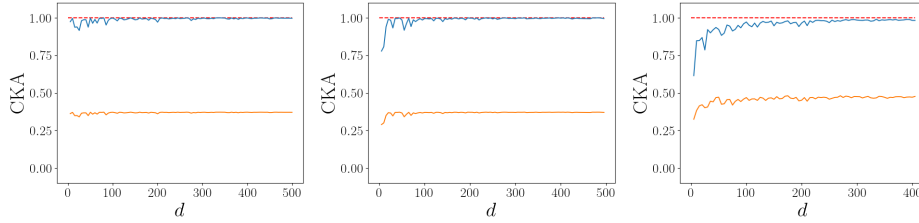


Figure 2: CKA between kernel matrices consisting of entries  $\Psi_{ij}$  and  $k_d^{(t)}(X_i, X_j)$  (Blue) and squared exponential kernel for control and  $k_d^{(t)}(X_i, X_j)$  (Orange) for three choices of  $A$  and  $R$ . (Left) Gaussian exponential family,  $A(\eta) = \eta^2/2$  and  $R(\eta) = \eta$ . (Middle) Bernoulli exponential family,  $A(\eta) = \log(1 + \exp \eta)$  and  $R(\eta) = \eta$ . (Right) Rectified Gaussian exponential family,  $A(\eta) = \eta^2$  and  $R(\eta) = \text{ReLU}(\eta)$ .

### 4.2 Inference using the DEK

We use the DEK for kernel ridge regression (Saunders, 1998) (KRR) (*cf* Gaussian process regression (Rasmussen & Williams, 2006)) on a suite of benchmarks. For each dataset, we first partition the data into an 80 – 20 train-test split. Using the training set, we perform 5-fold cross-validation for hyperparameter selection using the default settings of sci-kit learn’s `GridSearchCV`, which performs model selection based on the coefficient of determination. The hyperparameter grid we search over is described in Table 4, Appendix J.2. We then compute the RMSE on the held out test set using all training data. We repeat this procedure for 100 different random shuffles of dataset, and find the sample average and standard deviation RMSE over the random shuffles. The results are reported in Table 2. The input  $X$  is preprocessed by subtracting the sample average and dividing by the sample standard deviation of each feature. Additionally, the target data  $y$  is mean-centered and scaled by the sample standard deviation. The reported RMSE is after conversion of  $y$  back to original units.

Data	DEK or $\ell$ DEK	NNK	NTK	SEK	
yacht	<b><math>0.65 \pm 0.21</math></b>	$2.13 \pm 0.57$	$2.75 \pm 0.58$	$3.62 \pm 0.67$	0
diabetes	<b><math>54.51 \pm 3.29</math></b>	<b><math>54.58 \pm 3.30</math></b>	<b><math>55.05 \pm 3.38</math></b>	<b><math>54.75 \pm 3.32</math></b>	70
energy1	<b><math>1.00 \pm 0.11</math></b>	<b><math>1.01 \pm 0.11</math></b>	$1.67 \pm 0.14$	<b><math>1.08 \pm 0.13</math></b>	78
energy2	<b><math>1.58 \pm 0.15</math></b>	<b><math>1.58 \pm 0.15</math></b>	$2.10 \pm 0.18$	<b><math>1.58 \pm 0.16</math></b>	68
concrete	<b><math>4.94 \pm 0.47</math></b>	<b><math>4.97 \pm 0.47</math></b>	<b><math>5.05 \pm 0.48</math></b>	$5.65 \pm 0.39$	60
wine	<b><math>0.57 \pm 0.02</math></b>	$0.61 \pm 0.02$	<b><math>0.54 \pm 0.02</math></b>	$0.62 \pm 0.02$	1

Table 2: RMSE of KRR models ( $\pm$  one standard deviation over 100 random seeds). We use the DEK described in § G.2.5, which outperforms other kernels according to the sample average of the RMSE, although often the difference in performance is small compared with the standard deviation over 100 seeds. The final column is the number of times the best DEK found using `GridSearchCV` was an NNK.

Since the DEK is a strict generalisation of the NNK, we expect the DEK to strictly out-perform the NNK. We find that `GridSearchCV` sometimes picks out settings that correspond with an NNK, but often does not. The number of times `GridSearchCV` collapses the DEK to the NNK is indicated in the last column of Table 2. Our results are consistent with the previously established observation that “NNKs frequently outperform NTKs” (Lee et al., 2020). More interestingly, we find that for each dataset, the DEK performs as well or better than every other kernel, including the NNK.

## 5 Conclusion

We introduced the DEK, a kernel analogue of implicit neural network models. The DEK is defined as the limiting inner product between two features computed using a feature update procedure as the dimensionality of the features goes to infinity.

We considered the problem of whether a deterministic update procedure for the DEK exists (3), and whether this update rule converges (4). We focused on the special case where the features are latent variables in an exponential family PCA model (with not necessarily canonical link function) learnt using SGD. Leveraging the connection between infinitely wide explicit neural networks and kernel methods, we showed how in such a setting an explicit update rule can be computed. The update rule is a composition of functions involving NNK building blocks.

The DEK has a number of interesting properties. The DEK is able to recover instances of the NNK, and also resembles the NTK. Importantly, unlike the NNK and NTK, the deep layer structure of the DEK is motivated entirely from an optimisation perspective. The activation functions (and thus kernels) involved in the computation of the DEK can be related back to statistical modelling assumptions on the data through the exponential family. In particular, the activation functions share a connection to the log partition function and inverse link function of the exponential family. On a series of benchmarks, the DEK performs as well as or outperforms the NNK, NTK and SEK.

Our work admits several natural extensions. The matrix  $V$  which represents a linear transformation or fully connected layer may be constrained to resemble a convolutional layer, and we expect a convolutional variant of the DEK to be tractable (Novak et al., 2019). Since our construction is naturally probabilistic, the Laplace approximation about the MAP may yield a tractable means of obtaining principled uncertainty estimates for kernel methods beyond the regular Gaussian process framework. Since the DEK satisfies a fixed point equation, implicit differentiation may be used to compute derivatives of the DEK with respect to its hyperparameters, mirroring the neural network counterpart (Bai et al., 2019).

We hope that our optimisation view and deterministic kernel update rule stimulates new research in both deep learning and kernel methods.

## References

- Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*, 2020.
- James H Albert and Siddhartha Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American statistical Association*, 88(422):669–679, 1993.
- Luís B. Almeida. A learning rule for asynchronous perceptrons with feedback in a combinatorial environment. 1990.
- Sanjeev Arora, Simon S Du, Zhiyuan Li, Ruslan Salakhutdinov, Ruosong Wang, and Dingli Yu. Harnessing the power of infinitely wide deep nets on small-data tasks. *International Conference on Learning Representations*, 2020.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. Deep equilibrium models. *Advances in Neural Information Processing Systems*, 32, 2019.
- Arindam Banerjee, Srujana Merugu, Inderjit S Dhillon, Joydeep Ghosh, and John Lafferty. Clustering with bregman divergences. *Journal of machine learning research*, 6(10), 2005.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. *Advances in Neural Information Processing Systems*, 32, 2019.
- Youngmin Cho and Lawrence K. Saul. Kernel methods for deep learning. In *Advances in Neural Information Processing Systems*, pp. 342–350. 2009.
- Michael Collins, Sanjoy Dasgupta, and Robert E Schapire. A generalization of principal components analysis to the exponential family. *Advances in neural information processing systems*, 14, 2001.
- Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. Algorithms for learning kernels based on centered alignment. *The Journal of Machine Learning Research*, 13:795–828, 2012.
- Nello Cristianini, John Shawe-Taylor, Andre Elisseeff, and Jaz Kandola. On kernel-target alignment. *Advances in neural information processing systems*, 14, 2001.
- Ricky Der and Daniel Lee. Beyond gaussian processes: On the distributions of infinite networks. In *Advances in Neural Information Processing Systems*, volume 18. MIT Press, 2005.
- David Duvenaud. *Automatic model construction with Gaussian processes*. PhD thesis, University of Cambridge, 2014.
- Bradley Efron. Prediction, estimation, and attribution. *International Statistical Review*, 88:S28–S59, 2020.
- Bradley Efron and Trevor Hastie. *Computer Age Statistical Inference, Student Edition: Algorithms, Evidence, and Data Science*, volume 6. Cambridge University Press, 2021.
- Laurent El Ghaoui, Fangda Gu, Bertrand Travacca, Armin Askari, and Alicia Tsai. Implicit deep learning. *SIAM Journal on Mathematics of Data Science*, 3(3):930–958, 2021.
- Stefano Favaro, Sandra Fortini, and Stefano Peluchetti. Deep stable neural networks: large-width asymptotics and convergence rates. *arXiv preprint arXiv:2108.02316*, 2021.
- Stefano Favaro, Sandra Fortini, and Stefano Peluchetti. Neural tangent kernel analysis of shallow alpha-stable relu neural networks. *arXiv preprint arXiv:2206.08065*, 2022.

- Adrià Garriga-Alonso, Carl Edward Rasmussen, and Laurence Aitchison. Deep convolutional networks as shallow gaussian processes. In *International Conference on Learning Representations*, 2018.
- Stephen Gould, Richard Hartley, and Dylan Campbell. Deep declarative networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(8):3988–4004, 2021.
- Insu Han, Amir Zandieh, Jaehoon Lee, Roman Novak, Lechao Xiao, and Amin Karbasi. Fast neural kernel embeddings for general activations. *Advances in neural information processing systems*, 2022.
- James W Hardin and Joseph W M Hilbe. *Generalized Linear Models and Extensions*. College Station Tex: StataCorp Press, 4th edition, 2018.
- Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–8580, 2018.
- Motonobu Kanagawa, Philipp Hennig, Dino Sejdinovic, and Bharath K Sriperumbudur. Gaussian processes and kernel methods: A review on connections and equivalences. *arXiv preprint arXiv:1807.02582*, 2018.
- C. T. Kelley. *Iterative Methods for Linear and Nonlinear Equations*. Society for Industrial and Applied Mathematics, 1995. doi: 10.1137/1.9781611970944.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. 2014.
- Nicolas Le Roux and Yoshua Bengio. Continuous neural networks. In *Artificial Intelligence and Statistics*, pp. 404–411, 2007.
- Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes. In *International Conference on Learning Representations*, 2018.
- Jaehoon Lee, Samuel Schoenholz, Jeffrey Pennington, Ben Adlam, Lechao Xiao, Roman Novak, and Jascha Sohl-Dickstein. Finite versus infinite neural networks: an empirical study. *Advances in Neural Information Processing Systems*, 33:15156–15172, 2020.
- David JC MacKay. *Introduction to Gaussian processes*. 1998.
- Julien Mairal, Piotr Koniusz, Zaid Harchaoui, and Cordelia Schmid. Convolutional kernel networks. *Advances in neural information processing systems*, 27, 2014.
- Alexander G de G Matthews, Mark Rowland, Jiri Hron, Richard E Turner, and Zoubin Ghahramani. Gaussian process behaviour in wide deep neural networks. In *The International Conference on Learning Representations*, 2018.
- P. McCullagh and J.A. Nelder. *Generalized Linear Models, Second Edition*. Monographs on Statistics and Applied Probability Series. Chapman & Hall, 1989.
- Lassi Meronen, Christabella Irwanto, and Arno Solin. Stationary activations for uncertainty calibration in deep learning. *Advances in Neural Information Processing Systems*, 33:2338–2350, 2020.
- Shakir Mohamed, Zoubin Ghahramani, and Katherine A Heller. Bayesian exponential family pca. *Advances in neural information processing systems*, 21, 2008.
- Radford M Neal. *Bayesian learning for neural networks*. PhD thesis, University of Toronto, 1995.
- Frank Nielsen and Vincent Garcia. Statistical exponential families: A digest with flash cards. *arXiv preprint arXiv:0911.4863*, 2009.



- Roman Novak, Lechao Xiao, Jaehoon Lee, Yasaman Bahri, Greg Yang, Jiri Hron, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. Bayesian deep convolutional networks with many channels are Gaussian processes. In *The International Conference on Learning Representations*, 2019.
- OpenAI, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Jozefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique Ponde de Oliveira Pinto, Jonathan Raiman, Tim Salimans, Jeremy Schlatter, Jonas Schneider, Szymon Sidor, Ilya Sutskever, Jie Tang, Filip Wolski, and Susan Zhang. Dota 2 with large scale deep reinforcement learning. 2019.
- Tim Pearce, Russell Tsuchida, Mohamed Zaki, Alexandra Brintrup, and Andy Neely. Expressive priors in Bayesian neural networks: Kernel combinations and periodic functions. In *Uncertainty in Artificial Intelligence*, 2019.
- Stefano Peluchetti, Stefano Favaro, and Sandra Fortini. Stable behaviour of infinitely wide deep neural networks. In *International Conference on Artificial Intelligence and Statistics*, pp. 1137–1146. PMLR, 2020.
- Fernando Pineda. Generalization of back propagation to recurrent and higher order neural networks. In *Neural information processing systems*, 1987.
- R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2021. URL <https://www.R-project.org/>.
- Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *international conference on machine learning*, pp. 2847–2854. PMLR, 2017.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in neural information processing systems*, 20, 2007.
- Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine learning*. Adaptive computation and machine learning. MIT Press, 2006.
- Max Revay, Ruigang Wang, and Ian R Manchester. Lipschitz bounded equilibrium networks. *arXiv preprint arXiv:2010.01732*, 2020.
- Danilo Riccio, Matthias J Ehrhardt, and Martin Benning. Regularization of inverse problems: Deep equilibrium models versus bilevel learning. *arXiv preprint arXiv:2206.13193*, 2022.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.
- Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence*, 1(5):206–215, 2019.
- C Saunders. Ridge regression learning algorithm in dual variables. In *Proceedings of the 15th International Conference on Machine Learning*, 1998.
- Bernhard Schölkopf, Alexander J Smola, Francis Bach, et al. *Learning with kernels: support vector machines, regularization, optimization, and beyond*. MIT press, 2002.
- William Fleetwood Sheppard. On the application of the theory of error to cases of normal distribution and normal correlation. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character*, (192):140, 1899.
- Jascha Sohl-Dickstein, Roman Novak, Samuel S Schoenholz, and Jaehoon Lee. On the infinite width limit of neural networks with a standard parameterization. *arXiv preprint arXiv:2001.07301*, 2020.

- Peter Sollich. Bayesian methods for support vector machines: Evidence and predictive class probabilities. *Machine learning*, 46(1):21–52, 2002.
- Russell Tsuchida. *Results on infinitely wide multi-layer perceptrons*. PhD thesis, The University of Queensland, 2020.
- Russell Tsuchida and Cheng Soon Ong. Deep equilibrium models as estimators for continuous latent variables. In *arXiv preprint*, 2022.
- Russell Tsuchida, Fred Roosta, and Marcus Gallagher. Invariance of weight distributions in rectified MLPs. In *International Conference on Machine Learning*, pp. 5002–5011, 2018.
- Russell Tsuchida, Tim Pearce, Chris van der Heide, Fred Roosta, and Marcus Gallagher. Avoiding kernel fixed points: Computing with elu and gelu infinite networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 9967–9977, 2021.
- Russell Tsuchida, Suk Yee Yong, Mohammad Ali Armin, Lars Petersson, and Cheng Soon Ong. Declarative nets that are equilibrium models. In *International Conference on Learning Representations*, 2022.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2019. doi: 10.1017/9781108627771.
- Martin J Wainwright, Michael I Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008.
- Paul FV Wiemann, Thomas Kneib, and Julien Hambuckers. Using the softplus function to construct alternative link functions in generalized linear models and beyond. *arXiv preprint arXiv:2111.14207*, 2021.
- Christopher KI Williams. Computing with infinite networks. In *Advances in neural information processing systems*, pp. 295–301, 1997.
- Ezra Winston and J Zico Kolter. Monotone operator equilibrium networks. *Advances in Neural Information Processing Systems*, 33:10718–10728, 2020.
- Blake Woodworth, Suriya Gunasekar, Jason D Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pp. 3635–3673. PMLR, 2020.
- Stephen J. Wright and Benjamin Recht. *Optimization for Data Analysis*. Cambridge University Press, 2022. doi: 10.1017/9781009004282.
- Xingyu Xie, Qiuhaio Wang, Zenan Ling, Xia Li, Yisen Wang, Guangcan Liu, and Zhouchen Lin. Optimization induced equilibrium networks. *arXiv preprint arXiv:2105.13228*, 2021.
- Greg Yang. Wide feedforward or recurrent neural networks of any architecture are Gaussian processes. In *Advances in Neural Information Processing Systems*, pp. 9947–9960. 2019a.
- Greg Yang. Scaling limits of wide neural networks with weight sharing: Gaussian process behavior, gradient independence, and neural tangent kernel derivation. In *arXiv preprint arXiv:1902.04760*, 2019b.

Amir Zandieh, Insu Han, Haim Avron, Neta Shoham, Chaewon Kim, and Jinwoo Shin. Scaling neural tangent kernels via sketching and random features. *Advances in Neural Information Processing Systems*, 34:1062–1073, 2021.

## A Table of notations

Symbol	Name	Description
Features and feature updates		
$\psi_{X_i}$	Implicit function of $X_i$	For example, the solution to an optimisation problem, differential equation or root finding problem depending on $X_i$ .
$\psi_{X_i}^{(t+1)}$	Feature for input $X_i$	As in (1). The result of applying $t$ iterations of a numerical procedure that compute the implicit function $\psi_{X_i}$ .
$g^{(t)}$	Feature update rule	As in (1)
$\Psi$	Feature space	Features are an element of this space.
Kernels and kernel updates (component-wise)		
$\Psi_{ij}$	$\ell$ DEK evaluated at $X_i$ and $X_j$	As in (2).
$\Psi_{ij}^{(t+1)}$	DEK evaluated at $X_i$ and $X_j$	As in (2).
$G_{ij}$	DEK update rule (component-wise)	As in Theorem 4.
$\Psi$	DEK evaluation space	Evaluations of the DEK are an element of this space.
Kernels and kernel updates (matrix)		
$\Psi$	$\ell$ DEK matrix evaluated at $(X_1, X_2)$	$\Psi \in \mathbb{S}_+^2$ and the $ij$ th element of $\Psi$ is $\Psi_{ij}$ .
$\Psi^{(t+1)}$	DEK matrix evaluated at $(X_1, X_2)$	$\Psi^{(t+1)} \in \mathbb{S}_+^2$ and the $ij$ th element of $\Psi^{(t+1)}$ is $\Psi_{ij}^{(t+1)}$ .
$G$	DEK update rule (matrix version)	As in (3).
$\mathbb{S}_+^2$	Convex cone of PSD matrices	$\mathbb{S}_+^2 = \{M \in \mathbb{R}^{2 \times 2} \mid M \succeq 0\}$ .
Inputs and data		
$\mathbb{X}$	Input space	A space $\mathbb{X} \subseteq \mathbb{R}^l$ to which input belongs.
$X_i$	Input vector	An element of $\mathbb{X}$ .
$X$	Input matrix	An $N \times l$ matrix, where each row represents a single element of $\mathbb{X}$ .
$\Gamma$	Random mapping	$\Gamma : \mathbb{X} \rightarrow \mathbb{Y}$ is a random mapping which translates input $X$ to data $Y$ belonging to the support $\mathbb{Y}$ of an exponential family.
$Y = \Gamma(X)$	Data	The result of applying $\Gamma$ to input $X$
$\gamma_i(X)$	Random mapping coordinate evaluation	The $i$ th coordinate of $\Gamma(X)$ .
Exponential families		
$\mathbb{Y}$	Exponential family support	As in § 2.2. The space over which the exponential family distribution has non-zero mass. We take $\mathbb{Y} \subseteq \mathbb{R}$ .
$T$	Sufficient statistic	As in § 2.2.
$\eta$	Canonical parameter	As in § 2.2. The canonical parameter belongs to an open set $\mathcal{H} \subseteq \mathbb{R}$ .
$A$	Log partition function	The function that returns logarithm of the normalising constant of the exponential family as a function of its canonical parameter. $A : \mathbb{H} \rightarrow \mathbb{R}$ . As in § 2.2.
$s^{-1}$	Inverse link function	The function that maps a parameter of the exponential family to the expectation parameter of the exponential family. As in (8).

$R$	Canonical nonlinearity	As in Proposition 2. $R : \mathbb{R} \rightarrow \mathbb{H}$ maps the result of a linear transformation to a canonical parameter of the exponential family.
Activation functions		
$\zeta, \zeta_1, \zeta_2$	General activation function	A generic activation function of a neural network.
$\rho$	factor activation	The derivative of the canonical nonlinearity. That is, $\rho(a) = R'(a)$ .
$\sigma$	chain activation	The derivative of the composition of the log partition function and the canonical nonlinearity. That is, $\sigma(a) = (A \circ R)'(a)$
$u$	Heaviside step function	A special case of $\zeta$ . $u(a)$ takes the value of 0 if $a < 0$ , 1 if $a > 0$ , and 0.5 if $a = 0$ .
ReLU	Rectified linear unit	A special case of $\zeta$ . $\text{ReLU}(a) = u(a)a$ .
erf	Error function	$\text{erf}(a) = \frac{2}{\sqrt{\pi}} \int_0^a e^{-z^2} dz$ . As in § 3.3.

Table 3: Summary of notation used throughout this paper.

## B What is the relationship between $s^{-1}$ , $R$ , $\sigma$ and $\rho$ ?

**Nonlinearly parameterised exponential families are densities** Any member of a given exponential family is a density (mass) function. That is, for every  $\nu \in \mathbb{H}$ ,

$$\int p(y \mid \nu) dy = 1 \quad \text{and} \quad p(y \mid \nu) \geq 0. \quad (21)$$

Given a nonlinearity  $R : \mathbb{R} \rightarrow \mathbb{H}$ , it is immediate that any member of a given nonlinearly parameterised exponential family is a density. That is, for any  $a \in \mathbb{R}$ , defining  $\nu_0 = R(a) \in \mathbb{H}$ , from (21) we have

$$\int p(y \mid R(a)) dy = \int p(y \mid \underbrace{\nu_0}_{\in \mathbb{H}}) dy = 1 \quad \text{and} \quad p(y \mid R(a)) = p(y \mid \underbrace{\nu_0}_{\in \mathbb{H}}) \geq 0.$$

**Identities relating  $s^{-1}$ ,  $R$ ,  $\sigma$ , and  $\rho$**  The canonical link function is one which expresses the canonical parameter  $\eta$  in terms of the expectation parameter  $\mathbb{E}[T(Y) \mid \eta]$ . When  $R$  is the identity, we have that  $A'(\eta) = \mathbb{E}[T(Y) \mid \eta]$  and so the canonical link function is  $(A')^{-1}$ . That  $A'(\eta) = \mathbb{E}[T(Y) \mid \eta]$  follows from the fact that  $A$  is a cumulant generating function for the sufficient statistic (Wainwright et al., 2008, Proposition 3.1).

A (not necessarily canonical) link function is one which expresses a (not necessarily canonical) parameter  $a$  in terms of the expectation parameter  $\mathbb{E}[T(Y)]$ . We now discuss how given an exponential family and a link function can be related to a choice of  $R$ .

In the general setting, since  $A$  is a cumulant generating function, the inverse link function  $s^{-1}$  satisfies

$$A'(R(a)) = s^{-1}(a) \quad (22)$$

Noting that  $A'$  is invertible because  $A$  is strictly convex, (22) implies that for a desired link function  $s$ , we must choose

$$\begin{aligned} R(a) &= ((A')^{-1} \circ s^{-1})(a) \\ &= (A')^{-1}(\eta), \quad \eta = s^{-1}(a) \end{aligned} \quad (23)$$

Since  $\rho(a) = R'(a)$ , we have that

$$\begin{aligned} \rho(a) &= \frac{dR}{d\eta} \frac{d\eta}{da} \\ &= ((A')^{-1})'(s^{-1}(a)) (s^{-1})'(a) \\ &= \frac{(s^{-1})'(a)}{A'' \circ (A')^{-1} \circ s^{-1}(a)} \quad (\text{by the inverse function theorem}). \end{aligned} \quad (24)$$

Since  $\sigma(a) = (A \circ R)'(a)$  and  $\rho(a) = R'(a)$ , we have that

$$\begin{aligned} \sigma(a) &= \underbrace{A'(R(a))}_{s^{-1}(a)} \rho(a) \\ &= \frac{s^{-1}(a) (s^{-1})'(a)}{A'' \circ (A')^{-1} \circ s^{-1}(a)}. \end{aligned} \quad (25)$$

As expected, when  $s$  a canonical link function (which is to say that  $s^{-1}(a) = A'(a)$ ), (23) implies that  $R$  is the identity,  $\rho$  takes a constant value of 1 and  $\sigma$  is  $A'$ .

Some examples are given in Table 1.

## C Detailed neural network kernel description

Let  $W^{(1)} \in \mathbb{R}^{d \times n}$  be the weights of a fully connected hidden layer with activation function  $\sigma$ . Suppose each entry of  $W^{(1)}$  is i.i.d. with distribution  $\mathcal{N}(0, 1)$ . Given an input  $\phi_1 \in \mathbb{R}^{n \times 1}$  (we take the convention that vectors are column vectors), the signal in the hidden layer is  $h_1^{(1)} \triangleq \sigma(W^{(1)}\phi_1)$ . Given any two input features  $\phi_1$  and  $\phi_2$ , a normalised inner product of the features in the hidden layer is

$$\frac{1}{d} h_1^{(1)\top} h_2^{(1)} = \frac{1}{d} \sigma(W^{(1)}\phi_1)^\top \sigma(W^{(1)}\phi_2) = \frac{1}{d} \sum_{i=1}^d \sigma(W_i^\top \phi_1) \sigma(W_i^\top \phi_2), \quad (26)$$

where  $W_i^\top$  is the  $i$ th row of  $W^{(1)}$ . Note that since each row of  $W^{(1)}$  is i.i.d., (26) is an average of i.i.d. random variables. A strong law of large numbers says that the average of a sequence of i.i.d. random variables converges almost surely to the expectation if the expectation is finite. We therefore have that (26) converges almost surely to

$$k_\sigma(\phi_1, \phi_2) \triangleq \mathbb{E}_W [\sigma(W^\top \phi_1) \sigma(W^\top \phi_2)], \quad (27)$$

as  $d \rightarrow \infty$ . Here  $W^\top \in \mathbb{R}^{1 \times n}$  is a vector with i.i.d. entries drawn from  $\mathcal{N}(0, 1)$ . We call  $k_\sigma$  a single hidden layer neural network kernel (NNK) with activation function  $\sigma$ .

Note that  $W^\top$  is a row vector, and therefore  $W^\top \phi_1$  is a scalar. This means that while (27) is written as an expectation over  $n$ -variate random vector  $W$ , it is actually only an expectation over the bivariate random vector  $(\chi, \chi') = (W^\top \phi_1, W^\top \phi_2)$ . Since Gaussian random vectors are closed under affine transformations,  $(\chi, \chi')$  is a Gaussian random vector. The mean of each component is zero. The 2-by-2 covariance matrix  $\Sigma^{(1)}$  has entries

$$\Sigma_{12}^{(1)} = \mathbb{E}[(W^\top \phi_1)(W^\top \phi_2)] = \mathbb{E}\left[\sum_{p=1}^n \sum_{q=1}^n W_p \phi_{1p} W_q \phi_{2q}\right] = \sum_{p=1}^n \sum_{q=1}^n \mathbb{E}[W_p \phi_{1p} W_q \phi_{2q}] = \sum_{p=1}^n \mathbb{E}[W_p^2] \phi_{1p} \phi_{2p},$$

where the last equality is due to the fact that  $W_p$  and  $W_q$  are independent when  $p \neq q$ , and  $\phi_{1p}$  and  $\phi_{2q}$  are not random variables. Since  $\mathbb{E}W_p^2 = 1$ , the right most term is  $\phi_1^\top \phi_2$ . We may repeat a similar procedure for  $\Sigma_{11}$  and  $\Sigma_{22}$ , giving us an expression for the covariance

$$\Sigma^{(1)} \triangleq \begin{pmatrix} \Sigma_{11}^{(1)} & \Sigma_{12}^{(1)} \\ \Sigma_{12}^{(1)} & \Sigma_{22}^{(1)} \end{pmatrix} = \begin{pmatrix} \phi_1^\top \phi_1 & \phi_1^\top \phi_2 \\ \phi_2^\top \phi_1 & \phi_2^\top \phi_2 \end{pmatrix}. \quad (28)$$

It is instructive to rewrite (27) in two other forms. The first form explicitly shows the expectation with respect to the bivariate Gaussian which has covariance given by (28),

$$k_\sigma(\phi_1, \phi_2) = \mathbb{E}_{(\chi, \chi')^\top \sim \mathcal{N}(\mathbf{0}, \Sigma^{(1)})} [\sigma(\chi) \sigma(\chi')], \quad (29)$$

For the second form, we use notation to remind us that the kernel  $k_\sigma$  actually depends only on three scalar values. From (28), we observe that  $k_\sigma(\phi_1, \phi_2)$  depends on  $\phi_1$  and  $\phi_2$  *only* through the pairwise inner products  $\Sigma_{12}^{(1)} = \phi_1^\top \phi_2$ . Observe that by symmetry  $\Sigma_{12}^{(1)} = \Sigma_{21}^{(1)}$ . In other words, there exists a function  $\kappa_\sigma$  such that

$$k_\sigma(\phi_1, \phi_2) = \kappa_\sigma(\Sigma_{11}^{(1)}, \Sigma_{22}^{(1)}, \Sigma_{12}^{(1)}). \quad (30)$$

In summary, there are three equivalent ways to write an NNK,  $k_\sigma(\phi_1, \phi_2)$ :

- As an expectation over random vectors  $W$  corresponding to neural network weights (27),
- As an expectation over a bivariate Gaussian with covariance  $\Sigma^{(1)}$  (28),
- As a function of three arguments, explicitly showing the three parameters in the covariance of the bivariate Gaussian (30).

Closed-form expressions of  $k_\sigma$  for different  $\sigma$  are available (Williams, 1997; Le Roux & Bengio, 2007; Cho & Saul, 2009; Tsuchida et al., 2018; Pearce et al., 2019; Tsuchida, 2020; Meronen et al., 2020; Tsuchida et al., 2021; Han et al., 2022). For example, when  $\sigma$  is the ReLU function, the resulting kernel is known as the arc-cosine kernel of order 1 and is given by (Cho & Saul, 2009)

$$k_{\text{ReLU}}(\phi_1, \phi_2) = \frac{\|\phi_1\| \|\phi_2\|}{2\pi} (\sin \theta - (\pi - \theta) \cos \theta), \quad \text{where} \quad \cos \theta = \frac{\phi_1^\top \phi_2}{\|\phi_1\| \|\phi_2\|}.$$



## D Tools for concentration inequalities

The main purpose of this appendix is to introduce Bernstein's inequality and associated tools to apply to our problem at hand. We first need to introduce sub-Gaussian and sub-exponential random variables, and discuss special cases of how we may construct such random variables.

**Definition 9.** A centered random variable  $Y$  is sub-Gaussian if there exists an  $S > 0$  such that

$$\mathbb{E} \exp(Y^2/S^2) \leq 2.$$

The sub-Gaussian norm of  $Y$ ,

$$s \triangleq \inf \left\{ v > 0 : \mathbb{E} \exp(Y^2/v^2) \leq 2 \right\},$$

is the smallest  $S$ .

Bounded random variables are sub-Gaussian, and as an immediate consequence, so are constant random variables.

**Lemma 10** (Vershynin (2018) Example 2.5.8). A bounded random variable  $Y$  is sub-Gaussian with sub-Gaussian norm  $s$  satisfying

$$s \leq (\log 2)^{-1} \|Y\|_\infty,$$

where  $\|Y\|_\infty$  is the essential supremum of  $Y$ .

Lipschitz functions of Gaussian random variables are also sub-Gaussian.

**Lemma 11** (Vershynin (2018) Theorem 5.2.2). Let  $Y$  be a Gaussian random variable with variance  $a^2$ . Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be  $L$ -Lipschitz. Then  $f(Y) - \mathbb{E}f(Y)$  is sub-Gaussian with sub-Gaussian norm  $s_0$  satisfying

$$s_0 \leq C|a|L,$$

for some absolute constant  $C > 0$ . Furthermore, by the triangle inequality and Lemma 10  $f(Y)$  is sub-Gaussian and the sub-Gaussian norm  $s$  of  $f(Y)$  satisfies

$$s \leq C|a|L + (\log 2)^{-1} |\mathbb{E}f(Y)|.$$

A class of random variables which includes sub-Gaussian random variables is the class of sub-exponential random variables.

**Definition 12.** A random variable  $Y$  is sub-exponential if there exists an  $A > 0$  such that

$$\mathbb{E} \exp(|Y|/A) \leq 2$$

The sub-exponential norm of  $Y$ ,

$$a \triangleq \inf \left\{ v > 0 : \mathbb{E} \exp(|Y|/v) \leq 2 \right\},$$

is the smallest  $A$ .

Centering a sub-exponential random variable results in another sub-exponential random variable.

**Lemma 13** (Vershynin (2018) Exercise 2.7.10). If  $Y$  is sub-exponential with sub-exponential norm  $a_0$  then  $Y - \mathbb{E}Y$  is also sub-exponential, with sub-exponential norm  $a$  satisfying

$$a \leq Ca_0$$

for some absolute constant  $C > 0$ .

A useful fact is that a product of sub-Gaussian random variables is sub-exponential.

**Lemma 14** (Vershynin (2018) Lemma 2.7.7). *Let  $Y_1$  and  $Y_2$  be sub-Gaussian random variables with sub-Gaussian norms  $s_1$  and  $s_2$  respectively. Then their product  $Y_1 Y_2$  is sub-exponential with sub-exponential norm  $a$  satisfying  $a \leq s_1 s_2$ .*

Finally, sub-exponential random variables obey a useful quantitative form of a law of large numbers, which is a form of a Bernstein inequality.

**Theorem 15** (Bernstein's inequality, Corollary 2.8.3 of Vershynin (2018)). *Let  $Y_1, \dots, Y_d$  be a collection of random variables and write  $\mu_i = \mathbb{E}Y_i$  for  $i = 1, \dots, d$ . Suppose  $Y_1 - \mu_1, \dots, Y_d - \mu_d$  are independent sub-exponential random variables with sub-exponential norms  $a_1, \dots, a_d$ . Then, for every  $r \geq 0$ ,*

$$\mathbb{P}\left(\left|\frac{1}{d} \sum_{i=1}^d (Y_i - \mu_i)\right| \geq r\right) \leq 2 \exp\left(-cdM\right),$$

where  $M = \min\left\{\frac{r^2}{\max_i a_i^2}, \frac{r}{\max_i a_i}\right\}$  and  $c > 0$  is an absolute constant.

## E Analysis

The stochastic gradient  $\frac{\partial}{\partial \phi} L(\phi; X, \mathbf{V})$  evaluated at an arbitrary point  $\phi \in \Psi$  for input  $X$  and random  $\mathbf{V}$  is after scaling the sum of the gradient of the negative log prior and the stochastic gradient of the negative log likelihood,

$$\frac{\partial}{\partial \phi} L(\phi; X, \mathbf{V}) = \underbrace{\sqrt{\frac{m}{d}} \lambda \phi}_{\text{Gradient of negative log prior}} - \underbrace{\frac{1}{d} \mathbf{V}^\top (T(\Gamma(X)) \odot \rho(\mathbf{V}\phi) - \sigma(\mathbf{V}\phi))}_{\text{Stochastic gradient of log likelihood}}. \quad (31)$$

In order to prove Theorem 4, we will need to prove a series of lemmas. The intuition behind these lemmas is as follows. Assumption 2 means that if the limit were allowed to be applied, the gradient of the negative log prior term in (31) multiplied by the step size would look like  $\phi$ . This means that the update of SGD would just be the stochastic gradient of the log likelihood. We then examine the inner product of the stochastic gradient of the log likelihood, which would be the kernel update rule. The series of Lemmas is then as follows. We first convert the inner product of the stochastic gradient of the log likelihood to an approximate form that is easier to deal with (Lemma 16). We then confirm that the kernel update only involves the inner product of the stochastic gradients of the log likelihood (Lemma 17). Finally, we show that the inner products of the approximate form converges to a closed form update rule  $\mathbf{G}$  (Lemma 18). Assembling these lemmas together yields Theorem 4.

To this end, define the kernel

$$k_d(X_1, X_2; \phi_1, \phi_2) \triangleq \frac{1}{dm\lambda^2} (T(\Gamma(X_1)) \odot \rho(\mathbf{V}\phi_1) - \sigma(\mathbf{V}\phi_1))^\top \mathbf{V}\mathbf{V}^\top (T(\Gamma(X_2)) \odot \rho(\mathbf{V}\phi_2) - \sigma(\mathbf{V}\phi_2)),$$

which is a scaled inner product of the gradient of the negative log likelihood evaluated at inputs  $X_1$  and  $X_2$  and arbitrary points  $\phi_1$  and  $\phi_2$ . The factor  $\frac{1}{m} \mathbf{V}\mathbf{V}^\top \in \mathbb{R}^{d \times d}$  is approximately the identity matrix for large  $m$  under Assumption 1, leading to an easier to deal with approximation  $\tilde{k}_d(X_1, X_2; \phi_1, \phi_2)$  for  $k_d(X_1, X_2; \phi_1, \phi_2)$ ,

$$\tilde{k}_d(X_1, X_2; \phi_1, \phi_2) \triangleq \frac{1}{d\lambda^2} (T(\Gamma(X_1)) \odot \rho(\mathbf{V}\phi_1) - \sigma(\mathbf{V}\phi_1))^\top (T(\Gamma(X_2)) \odot \rho(\mathbf{V}\phi_2) - \sigma(\mathbf{V}\phi_2)).$$

Lemma 16 says that this approximation is exact in the infinite  $d$  limit, and quantifies the quality of this approximation when  $d$  is finite.

**Lemma 16.** *Let  $\phi_1 \in \mathbb{R}^m$  and  $\phi_2 \in \mathbb{R}^m$  be arbitrary and suppose Assumption 3 holds.*

**(16A)** *Under Assumption 1,  $k_d(X_1, X_2; \phi_1, \phi_2)$  converges in probability to  $\tilde{k}_d(X_1, X_2; \phi_1, \phi_2)$ .*

**(16B)** *Under Assumption 5, there exist constants  $Q > 0$  and  $c > 0$  such that for all  $\delta > 0$  and  $\epsilon_2 > 0$ ,*

$$\mathbb{P}\left(|k_d(X_1, X_2; \phi_1, \phi_2) - \tilde{k}_d(X_1, X_2; \phi_1, \phi_2)| \geq \frac{(K + \epsilon_2)}{\lambda^2} (2\epsilon_1 + \epsilon_1^2)\right) \leq 2 \exp\left(-cdM\right) + e^{-m\delta^2/2},$$

$$\text{where } \epsilon_1 = \sqrt{\frac{d}{m}} + \delta \text{ and } M = \min\left\{\frac{\epsilon_2^2}{Q^2}, \frac{\epsilon_2}{Q}\right\}.$$

*Proof.* We use the shorthand  $\Gamma_1 = \Gamma(X_1)$  and  $\Gamma_2 = \Gamma(X_2)$ . We have

$$\begin{aligned} & |k_d(X_1, X_2; \phi_1, \phi_2) - \tilde{k}_d(X_1, X_2; \phi_1, \phi_2)| \\ &= \frac{1}{d\lambda^2} |(T(\Gamma_1) \odot \rho(\mathbf{V}\phi_1) - \sigma(\mathbf{V}\phi_1))^\top \left(\frac{1}{m} \mathbf{V}\mathbf{V}^\top - \mathbf{I}\right) (T(\Gamma_2) \odot \rho(\mathbf{V}\phi_2) - \sigma(\mathbf{V}\phi_2))| \\ &\leq \frac{1}{d\lambda^2} \|T(\Gamma_1) \odot \rho(\mathbf{V}\phi_1) - \sigma(\mathbf{V}\phi_1)\| \left\| \frac{1}{m} \mathbf{V}\mathbf{V}^\top - \mathbf{I} \right\| \|T(\Gamma_2) \odot \rho(\mathbf{V}\phi_2) - \sigma(\mathbf{V}\phi_2)\| \\ &\leq \frac{1}{d\lambda^2} \max_{\phi_1, \phi_2} \left\{ \|T(\Gamma_1) \odot \rho(\mathbf{V}\phi_1) - \sigma(\mathbf{V}\phi_1)\|^2, \|T(\Gamma_2) \odot \rho(\mathbf{V}\phi_2) - \sigma(\mathbf{V}\phi_2)\|^2 \right\} \left\| \frac{1}{m} \mathbf{V}\mathbf{V}^\top - \mathbf{I} \right\| \\ &= \frac{1}{\lambda^2} \max_{\hat{\phi} \in \{\phi_1, \phi_2\}} \left( K_{\hat{\phi}} - \mathbb{E}[K_{\hat{\phi}}] \right) \left\| \frac{1}{m} \mathbf{V}\mathbf{V}^\top - \mathbf{I} \right\| + \mathbb{E}[K_{\hat{\phi}}] \left\| \frac{1}{m} \mathbf{V}\mathbf{V}^\top - \mathbf{I} \right\|, \end{aligned} \quad (32)$$

where  $K_{\hat{\phi}} = \frac{1}{d} \sum_{i=1}^d \left( T(\gamma_i(X_1)) \odot \rho(V_i^\top \hat{\phi}) - \sigma(V_i^\top \hat{\phi}) \right)^2$  and  $V_i^\top$  is the  $i$ th row of  $\mathbf{V}$ . The quantity  $\mathbb{E}[K_{\hat{\phi}}]$  is finite by Assumption 3.

Using a standard result (Wainwright, 2019, Example 6.2), we have that

$$\mathbb{P} \left( \left\| \mathbf{I} - \frac{1}{m} \mathbf{V} \mathbf{V}^\top \right\| \geq 2\epsilon_1 + \epsilon_1^2 \right) \leq e^{-m\delta^2/2}, \quad \epsilon_1 = \sqrt{\frac{d}{m}} + \delta. \quad (33)$$

Combining (32) and (33) and taking  $d \rightarrow \infty$  under Assumption 1, we have **(16A)**.

We may apply a Bernstein concentration inequality to  $K_{\hat{\phi}} - \mathbb{E}[K_{\hat{\phi}}]$  as follows. The variables  $T(\gamma_i(X_1))\rho(V_i^\top \psi) - \sigma(V_i^\top \psi)$  for each  $i$  are mutually independent. The quantities  $V_i^\top \hat{\psi}$  are zero-mean Gaussian (since Gaussian random variables are closed under linear combinations). By Assumption 5, each variable in the sum contains sub-Gaussian elements since bounded random variables are sub-Gaussian (Lemma 10), and Lipschitz functions of Gaussian random variables are sub-Gaussian (Lemma 11). The square of sub-Gaussian random variables is sub-exponential (Lemma 14). Sub-exponential random variables that are centered by subtracting their mean are also sub-exponential (Lemma 13). Therefore, by Bernstein's Theorem (Theorem 15), there exist constants  $c, Q > 0$  (depending on  $\rho, \sigma, X_1$  and  $X_2$ ) such that for every  $\epsilon_2 \geq 0$ ,

$$\mathbb{P} \left( \max_{\hat{\phi} \in \{\phi_1, \phi_2\}} |K_{\hat{\phi}} - \mathbb{E}K_{\hat{\phi}}| \geq \epsilon_2 \right) \leq 2 \exp(-cdM), \quad (34)$$

where  $M = \min \left\{ \frac{\epsilon_2^2}{Q^2}, \frac{\epsilon_2}{Q} \right\}$ .

Finally, combining (32), (33) and (34) via a union bound, we have

$$\mathbb{P} \left( |k_d(X_1, X_2; \phi_1, \phi_2) - \tilde{k}_d(X_1, X_2; \phi_1, \phi_2)| \leq \frac{(K + \epsilon_2)}{\lambda^2} (2\epsilon_1 + \epsilon_1^2) \right) \geq 1 - 2 \exp(-cdM) - e^{-m\delta^2/2}.$$

□

We now confirm that the kernel update rule only involves the inner product of the stochastic gradient of the log likelihood, and not the gradient of the log prior.

**Lemma 17.** *Suppose Assumptions 1, 2 (a) and 3 hold. Then applying SGD to objective (17),*

$$\Psi_{ij}^{(t+1)} = \lim_{d \rightarrow \infty} k_d(X_i, X_j; \psi_{X_i}^{(t)}, \psi_{X_j}^{(t)}) = \text{plim}_{d \rightarrow \infty} \tilde{k}_d(X_i, X_j; \psi_{X_i}^{(t)}, \psi_{X_j}^{(t)})$$

*Proof.* Combining the SGD update (15) and the stochastic gradient (31), we have that for input  $X$ , the  $t + 1$ th iterate of SGD satisfies

$$\begin{aligned} \psi_X^{(t+1)} &= \psi_X^{(t)} \left( 1 - \alpha^{(t)} \sqrt{\frac{m}{d}} \lambda \right) + \alpha^{(t)} \frac{1}{d} \mathbf{V}^{(t)\top} \left( T(\Gamma(X)) \odot \rho(\mathbf{V}^{(t)} \psi_X^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_X^{(t)}) \right) \\ \psi_X^{(t+1)} - \psi_X^{(t)} \left( 1 - \alpha^{(t)} \sqrt{\frac{m}{d}} \lambda \right) &= \alpha^{(t)} \frac{1}{d} \mathbf{V}^{(t)\top} \left( T(\Gamma(X)) \odot \rho(\mathbf{V}^{(t)} \psi_X^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_X^{(t)}) \right) \end{aligned} \quad (35)$$

Evaluating (35) at input  $X_i$  and  $X_j$ , and taking inner products, we find

$$\begin{aligned} & \left( \psi_{X_i}^{(t+1)} - \psi_{X_i}^{(t)} \left( 1 - \alpha^{(t)} \sqrt{\frac{m}{d}} \lambda \right) \right)^\top \left( \psi_{X_j}^{(t+1)} - \psi_{X_j}^{(t)} \left( 1 - \alpha^{(t)} \sqrt{\frac{m}{d}} \lambda \right) \right) \\ &= \alpha^{(t)2} \frac{1}{d^2} \left( T(\Gamma(X_i)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}) \right)^\top \mathbf{V}^{(t)} \mathbf{V}^{(t)\top} \left( T(\Gamma(X_j)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}) \right) \end{aligned}$$

Invoking Assumption 2, we see that the left hand side satisfies

$$\begin{aligned}
& \lim_{d \rightarrow \infty} \left( \psi_X^{(t+1)} - \psi_X^{(t)} (1 - \alpha^{(t)} \sqrt{\frac{m}{d}} \lambda) \right)^\top \left( \psi_{X_j}^{(t+1)} - \psi_{X_j}^{(t)} (1 - \alpha^{(t)} \sqrt{\frac{m}{d}} \lambda) \right) \\
&= \lim_{d \rightarrow \infty} \underbrace{\psi_X^{(t+1)\top} \psi_{X_j}^{(t+1)} + \psi_X^{(t)\top} \psi_{X_j}^{(t)} (1 - \alpha^{(t)} \sqrt{\frac{m}{d}} \lambda)^2}_{\rightarrow 0} \\
&\quad - \underbrace{(1 - \alpha^{(t)} \sqrt{\frac{m}{d}} \lambda)}_{\rightarrow 0} \left( \|\psi_X^{(t+1)}\| \|\psi_{X_j}^{(t)}\| a_1 + \|\psi_X^{(t)}\| \|\psi_{X_j}^{(t+1)}\| a_2 \right) \\
&= \Psi_{ij}^{(t+1)}
\end{aligned}$$

where  $a_1$  and  $a_2$  are cosine angles belonging to  $[-1, 1]$ . On the other hand, under Assumption 2 the right hand side satisfies

$$\begin{aligned}
& \lim_{d \rightarrow \infty} \alpha^{(t)^2} \frac{1}{d^2} (T(\Gamma(X_i)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}))^\top \mathbf{V}^{(t)} \mathbf{V}^{(t)\top} (T(\Gamma(X_j)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_j}^{(t)})) \\
&= \lim_{d \rightarrow \infty} \frac{1}{dm\lambda^2} (T(\Gamma(X_i)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}))^\top \mathbf{V}^{(t)} \mathbf{V}^{(t)\top} (T(\Gamma(X_j)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_j}^{(t)})) \\
&= \lim_{d \rightarrow \infty} k_d(X_i, X_j; \psi_{X_i}^{(t)}, \psi_{X_j}^{(t)}).
\end{aligned}$$

By Lemma 16, this limit is well defined and is given by  $\text{plim}_{d \rightarrow \infty} \tilde{k}_d(X_i, X_j; \psi_{X_i}^{(t)}, \psi_{X_j}^{(t)})$ .  $\square$

Finally, we show that the approximate form of inner products  $\tilde{k}_d$  converges to a closed form update rule  $\mathbf{G}$ .

**Lemma 18.** *Suppose Assumptions 1, 3 and 4 hold. Then*

$$\text{plim}_{d \rightarrow \infty} \tilde{k}_d(X_i, X_j; \psi_{X_i}^{(t)}, \psi_{X_j}^{(t)}) = G(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}; X_i, X_j),$$

where

$$\begin{aligned}
& G(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}; X_i, X_j) \\
&= \frac{1}{\lambda^2} \left( C_{ij} \kappa_\rho(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}) - \kappa_{\sigma, \rho}(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}) \mu_i - \kappa_{\rho, \sigma}(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}) \mu_j + \kappa_\sigma(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}) \right).
\end{aligned}$$

*Proof.* We have

$$\begin{aligned}
& \text{plim}_{d \rightarrow \infty} \tilde{k}_d(X_i, X_j; \psi_{X_i}^{(t)}, \psi_{X_j}^{(t)}) \\
&= \text{plim}_{d \rightarrow \infty} \frac{1}{d\lambda^2} (T(\Gamma(X_i)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}))^\top (T(\Gamma(X_j)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}))
\end{aligned}$$

Expanding the quadratic, we find

$$\begin{aligned}
\Psi_{ij}^{(t+1)} &= \text{plim}_{d \rightarrow \infty} \frac{1}{d\lambda^2} (T(\Gamma(X_i)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}))^\top (T(\Gamma(X_j)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}) - \sigma(\mathbf{V}^{(t)} \psi_{X_j}^{(t)})) \\
&= \text{plim}_{d \rightarrow \infty} \frac{1}{d\lambda^2} \left( (T(\Gamma(X_i)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}))^\top (T(\Gamma(X_j)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_j}^{(t)})) - \right. \\
&\quad (T(\Gamma(X_i)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}))^\top \sigma(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}) - \\
&\quad (T(\Gamma(X_j)) \odot \rho(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}))^\top \sigma(\mathbf{V}^{(t)} \psi_{X_i}^{(t)}) + \\
&\quad \left. \sigma(\mathbf{V}^{(t)} \psi_{X_i}^{(t)})^\top \sigma(\mathbf{V}^{(t)} \psi_{X_j}^{(t)}) \right)
\end{aligned}$$

The collection of  $d$  pairs  $\{(\mathbf{V}^{(t)}\psi_{X_i}^{(t)}, \mathbf{V}^{(t)}\psi_{X_j}^{(t)})_p\}_{p=1}^d$  is mutually independent and Gaussian given  $\psi_{X_i}^{(t)}$  and  $\psi_{X_j}^{(t)}$ . Letting  $V^\top \in \mathbb{R}^m$  be equal in distribution to a row of  $\mathbf{V}^{(t)}$ , a law of large numbers says that

$$\begin{aligned}\Psi_{ij}^{(t+1)} &= \frac{1}{\lambda^2} \left( c(X_i, X_j) \mathbb{E}_{V^\top} \left[ \lim_{d \rightarrow \infty} \rho(V^\top \psi_{X_i}^{(t)}) \rho(V^\top \psi_{X_j}^{(t)}) \right] - \right. \\ &\quad \mu(X_i) \mathbb{E}_{V^\top} \left[ \lim_{d \rightarrow \infty} \rho(V^\top \psi_{X_i}^{(t)}) \sigma(V^\top \psi_{X_j}^{(t)}) \right] - \\ &\quad \mu(X_j) \mathbb{E}_{V^\top} \left[ \lim_{d \rightarrow \infty} \rho(V^\top \psi_{X_j}^{(t)}) \sigma(V^\top \psi_{X_i}^{(t)}) \right] + \\ &\quad \left. \mathbb{E}_{V^\top} \left[ \lim_{d \rightarrow \infty} \sigma(V^\top \psi_{X_i}^{(t)}) \sigma(V^\top \psi_{X_j}^{(t)}) \right] \right).\end{aligned}$$

Now observe that conditional on  $\psi_{X_i}^{(t)}$  and  $\psi_{X_j}^{(t)}$ , the random vector  $(\chi, \chi')^\top = \lim_{d \rightarrow \infty} (V^\top \psi_{X_i}^{(t)}, V^\top \psi_{X_j}^{(t)})^\top$  is bivariate Gaussian with mean 0 and covariance matrix

$$\begin{pmatrix} \Psi_{ii}^{(t)} & \Psi_{ij}^{(t)} \\ \Psi_{ij}^{(t)} & \Psi_{jj}^{(t)} \end{pmatrix}.$$

Therefore, the terms on the right hand side involve evaluations of the functions  $\kappa_{\rho, \rho}, \kappa_{\rho, \sigma}, \kappa_{\sigma, \rho}$  and  $\kappa_{\sigma, \sigma}$ .  $\square$

Chaining Lemmas 17 and 18, we obtain our main theorem.

**Theorem 4.** *Suppose Assumptions 1, 2 (a), 3, and 4 hold. Let  $C_{ij} = c(X_i, X_j)$  and  $\mu_i = \mu(X_i)$  be as defined in Assumption 4. Then applying SGD to objective (17), the update rule  $\mathbf{G}$  (3) exists and can be decomposed into  $G$  (5) satisfying*

$$\begin{aligned}&G(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}; X_i, X_j) \\ &= \frac{1}{\lambda^2} \left( C_{ij} \kappa_{\rho}(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) - \kappa_{\sigma, \rho}(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) \mu_i - \kappa_{\rho, \sigma}(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) \mu_j + \kappa_{\sigma}(\Phi_{ii}, \Phi_{jj}, \Phi_{ij}) \right).\end{aligned}$$

Here  $\kappa_{\sigma}$ ,  $\kappa_{\rho}$ ,  $\kappa_{\sigma, \rho}$  and  $\kappa_{\rho, \sigma}$  are as defined by (11), (19) and Proposition 2.

*Proof.* By Lemma 17, under Assumptions 1, 2 (a) and 3, we have

$$\Psi_{12}^{(t+1)} = \lim_{d \rightarrow \infty} k_d(X_1, X_2; \psi_{X_1}^{(t)}, \psi_{X_2}^{(t)}) = \text{plim}_{d \rightarrow \infty} \tilde{k}_d(X_1, X_2; \psi_{X_1}^{(t)}, \psi_{X_2}^{(t)}).$$

By Lemma 18, under Assumptions 1, 3 and 4 we have

$$\text{plim}_{d \rightarrow \infty} \tilde{k}_d(X_1, X_2; \psi_{X_1}^{(t)}, \psi_{X_2}^{(t)}) = G(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}; X_i, X_j).$$

$\square$

The effect of the finite approximation is not described in Theorem 4. In order to describe the effect of finite approximations, we combine the previously proven Lemma 16B with the following Lemma 19, assembling them in Theorem 8.

We now turn to analysing the sensitivity of the fDEK when initialised around the  $\ell$ DEK. For this, we combine Lemma 16 with the following lemma.

**Lemma 19.** *Suppose Assumptions 1, 2 (a), 3, 4 and 5 hold. Then for all  $\epsilon > 0$  there exists some  $Q > 0$  (depending on  $\rho$ ,  $\sigma$ ,  $X_i$ , and  $X_j$ ) such that*

$$\mathbb{P}\left(\left|\tilde{k}_d(X_1, X_2; r, r') - \Psi_{12}\right| \geq \epsilon\right) \leq 2 \exp\left(-dcM\right),$$

where  $M = \min\left\{\frac{\epsilon^2}{Q^2}, \frac{\epsilon}{Q}\right\}$  and  $c > 0$  is some absolute constant.

*Proof.* We first write  $\tilde{k}_d$  as a sum. Letting  $V_i^\top \in \mathbb{R}^m$  denote the  $i$ th row of  $V$  and  $\gamma_i(X)$  denote the  $i$ th coordinate of  $\Gamma(X)$ ,

$$\tilde{k}_d(X_1, X_2; r, r') = \frac{1}{d\lambda^2} \sum_{i=1}^d (T(\gamma_i(X_1))\rho(V_i^\top r) - \sigma(V_i^\top r))(T(\gamma_i(X_2))\rho(V_i^\top r') - \sigma(V_i^\top r')).$$

In this form, we observe that  $\mathbb{E}\tilde{k}_d(X_1, X_2; r, r') = G(\Psi_{11}, \Psi_{22}, \Psi_{12}; X_1, X_2) = \Psi_{12}$ . We therefore seek to concentrate  $\tilde{k}_d(X_1, X_2; r, r')$  about its mean.

The bivariate pairs  $(T(\gamma_i(X_1))\rho(V_i^\top r) - \sigma(V_i^\top r), T(\gamma_i(X_2))\rho(V_i^\top r') - \sigma(V_i^\top r'))$  are independent from every other pair. The quantities  $V_i^\top r$  and  $V_i^\top r'$  are zero-mean Gaussian (since Gaussian random variables are closed under linear combinations). Each pair contains sub-Gaussian elements since bounded random variables are sub-Gaussian (Lemma 10), and Lipschitz functions of Gaussian random variables are sub-Gaussian (Lemma 11). The product of two sub-Gaussian random variables is sub-exponential (Lemma 14). Sub-exponential random variables that are centered by subtracting their mean are also sub-exponential (Lemma 13). Therefore, by Bernstein's Theorem (Theorem 15), there exist constants  $c, Q > 0$  (depending on  $\rho, \sigma, X_i$ , and  $X_j$ ) such that for every  $\epsilon \geq 0$ ,

$$\mathbb{P}\left(|\tilde{k}_d(X_1, X_2; r, r') - \Psi_{12}| \geq \epsilon\right) \leq 2 \exp\left(-cdM\right),$$

where  $M = \min\left\{\frac{\epsilon^2}{Q^2}, \frac{\epsilon}{Q}\right\}$ . □

**Theorem 8.** Suppose Assumptions 1, 2 (b), 3, 4 and 5 hold. Let initial guesses be  $\psi_{X_1}^{(0)} = r_1$  and  $\psi_{X_2}^{(0)} = r_2$  as in Definition 7. Then there exist constants  $Q_2, Q_3, c_2, c_3 > 0$  such that for all  $\delta > 0, \epsilon_2 > 0$  and  $\epsilon_2$ ,

$$\mathbb{P}\left(|\bar{\Psi}_{12}^{(1)} - \Psi_{12}| \leq \epsilon_1 + \epsilon_2\right) \geq 1 - \delta_1 - \delta_2,$$

where

$$\epsilon_1 = \frac{K + \epsilon_2}{\lambda^2} (2\epsilon_1 + \epsilon_1^2), \quad \delta_1 = 2 \exp\left(-c_2 d M_2\right) + \exp\left(-m\delta^2/2\right) \quad \text{and} \quad \delta_2 = 2 \exp\left(-dc_3 M_3\right)$$

and  $\epsilon_1 = \sqrt{\frac{d}{m}} + \delta$ ,  $M_2 = \min\left\{\frac{\epsilon_2^2}{Q_2^2}, \frac{\epsilon_2}{Q_2}\right\}$  and  $M_3 = \min\left\{\frac{\epsilon_2^2}{Q_3^2}, \frac{\epsilon_2}{Q_3}\right\}$  and  $c_3 > 0$  is some absolute constant.

*Proof.* Under Assumption 2 (b), plugging the stochastic gradient (31) into the SGD update rule (35), we have

$$\begin{aligned} \bar{\Psi}_{ij}^{(t+1)} &= \frac{1}{\lambda^2 dm} (T(\Gamma(X_i)) \odot \rho(V^{(t)}\psi_{X_i}^{(t)}) - \sigma(V^{(t)}\psi_{X_i}^{(t)}))^\top V^{(t)} V^{(t)\top} (T(\Gamma(X_j)) \odot \rho(V^{(t)}\psi_{X_j}^{(t)}) - \sigma(V^{(t)}\psi_{X_j}^{(t)})) \\ &= k_d(X_i, X_j; \psi_{X_i}^{(t)}, \psi_{X_j}^{(t)}) \\ \bar{\Psi}_{ij}^{(1)} &= k_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)}). \end{aligned}$$

By Lemma 16, we may approximate  $k_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)})$  by  $\tilde{k}_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)})$  with high probability. By Lemma 19, we may approximate  $\tilde{k}_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)})$  by  $\Psi_{ij}$  with high probability. The proof then follows by applying a triangle inequality and a union bound.

In more detail, the triangle inequality says

$$|\bar{\Psi}_{12}^{(1)} - \Psi_{12}| \leq |k_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)}) - \tilde{k}_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)})| + |\tilde{k}_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)}) - \Psi_{ij}|. \quad (36)$$

Let  $\delta > 0$  and  $\epsilon_2 > 0$  be arbitrary. Define  $\epsilon_1 = \sqrt{\frac{d}{m}} + \delta$ . Define  $\epsilon_1 = \frac{K + \epsilon_2}{\lambda^2} (2\epsilon_1 + \epsilon_1^2)$ . Let  $A_1$  denote the event that  $|k_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)}) - \tilde{k}_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)})| \geq \epsilon_1$ . Then by Lemma 16 there exists some constants  $Q_2 > 0$  and  $c_2 > 0$  such that

$$\mathbb{P}(A_1) \leq \delta_1, \quad (37)$$

where  $\delta_1 = 2 \exp(-c_2 d M_2) + \exp(-m \delta^2/2)$  and  $M_2 = \min\{\frac{\varepsilon_2^2}{Q_2^2}, \frac{\varepsilon_2}{Q_2}\}$ .

Let  $\varepsilon_2 > 0$  be arbitrary. Let  $A_2$  denote the event that  $|\tilde{k}_d(X_i, X_j; \psi_{X_i}^{(0)}, \psi_{X_j}^{(0)}) - \Psi_{ij}| \geq \varepsilon_2$ . Then by Lemma 19, there exists some  $Q_3 > 0$  such that

$$\mathbb{P}(A_2) \leq \delta_2, \tag{38}$$

where  $\delta_2 = 2 \exp(-dc_3 M_3)$ ,  $M_3 = \min\{\frac{\varepsilon_2^2}{Q_3^2}, \frac{\varepsilon_2}{Q_3}\}$  and  $c_3 > 0$  is an absolute constant.

Combining (37) and (38) by a union bound, we obtain

$$\begin{aligned} \mathbb{P}(A_1 \cup A_2) &\leq \delta_1 + \delta_2 \\ \mathbb{P}(A_1^c \cap A_2^c) &\geq 1 - \delta_1 - \delta_2 \end{aligned}$$

Finally, if  $A_1^c \cap A_2^c$  then by (36),  $|\bar{\Psi}_{12}^{(1)} - \Psi_{12}| \leq \varepsilon_1 + \varepsilon_2$ . □



## F Random finite forms for the explicit kernel $c$ and mean function $\mu$

Suppose the sufficient statistic  $T$  is the identity.

**Linear kernel.** Let  $\Gamma(X_1) = QX_1$ , where each entry of  $Q \in \mathbb{R}^{d \times l}$  is sampled i.i.d. from  $\mathcal{N}(0, v^2)$ . Then we obtain the linear kernel,

$$\begin{aligned} c(X_1, X_2) &= v^2 \lim_{d \rightarrow \infty} \frac{1}{d} X_1^\top Q^\top Q X_2 \\ &= v^2 X_1^\top X_2. \end{aligned}$$

In this case,  $\mu(X_1) = 0$ . since  $\frac{1}{d} \sum_{i=1}^d \mathbb{E}[Q_i^\top X_1] = \frac{1}{d} \mathbb{E}[Q_i^\top] X_1 = 0$ , where  $Q_i^\top$  is the  $i$ th row of  $Q$ .

**Squared exponential kernel.** We may obtain stationary nonlinear kernels via a random Fourier feature type construction (Rahimi & Recht, 2007). Suppose  $d$  is even and  $Q \in \mathbb{R}^{d/2 \times l}$  is sampled i.i.d. from  $\mathcal{N}(0, 1)$ . Define

$$\Gamma(X_1) = A \odot \begin{pmatrix} \cos(QX_1) \\ \sin(QX_1) \end{pmatrix} \in \mathbb{R}^{d \times 1},$$

where elements of  $A = (a_1, \dots, a_{d/2}, b_1, \dots, b_{d/2})^\top \in \mathbb{R}^{d \times 1}$  are sampled i.i.d. from  $\mathcal{N}(\mu_v, v^2)$ . Then  $\mathbb{E}[a_i^2] = v^2 + \mu_v^2$  and

$$\begin{aligned} c(X_1, X_2) &= \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^{d/2} a_i^2 \cos(Q_i^\top X_1) \cos(Q_i^\top X_2) + \frac{1}{d} \sum_{j=1}^{d/2} b_j^2 \sin(Q_j^\top X_1) \sin(Q_j^\top X_2) \\ &= \frac{v^2 + \mu_v^2}{2} \mathbb{E}[\cos(Q^\top X_1) \cos(Q^\top X_2) + \sin(Q^\top X_1) \sin(Q^\top X_2)], \quad Q \sim \mathcal{N}(0, I) \\ &= \frac{v^2 + \mu_v^2}{2} \mathbb{E}[\cos(Q^\top (X_1 - X_2))] \\ &= \frac{v^2 + \mu_v^2}{2} \exp\left(-\frac{1}{2} \|X_1 - X_2\|_2^2\right). \end{aligned}$$

An extension to arbitrary stationary kernels follows using Bochner's theorem to define the probability measure of  $Q$  via a Fourier transform (Rahimi & Recht, 2007). An extension to arbitrary covariance structures can be obtained by introducing a dependency structure among elements of rows of  $Q$ .

The reason that  $A$  is introduced is to allow the mean to converge to zero, so that  $\mu(X) = 0$  can be realised. That is,

$$\begin{aligned} \mu(X_1) &= \lim_{d \rightarrow \infty} \frac{1}{d} \sum_{i=1}^{d/2} a_i \cos(Q_i^\top X_1) + \frac{1}{d} \sum_{j=1}^{d/2} b_j \sin(Q_j^\top X_1) \\ &= \frac{\mu_v}{2} \mathbb{E}[\cos(Q^\top X_1)] \\ &= \frac{\mu_v}{2} \exp\left(-\frac{1}{2} \|X_1\|^2\right) \end{aligned}$$

is zero whenever  $\mu_v = 0$ .

### F.1 A model we found to be practically useful

Let  $\Gamma(X_1) = A \cdot \text{ReLU}(QX_1)$ , where elements of  $A$  are sampled i.i.d. from  $\mathcal{N}(\mu_v, v^2)$ . Then  $c$  is the arc-cosine kernel of degree 1 (Cho & Saul, 2009) (see (47)),

$$c(X_1, X_2) = (v^2 + \mu_v^2) \frac{\|X_1\| \|X_2\|}{2\pi} (\sin \theta + (\pi - \theta) \cos \theta), \quad \theta = \arccos \frac{X_1^\top X_2}{\|X_1\| \|X_2\|}.$$

The mean function is given by

$$\begin{aligned}\mu(X_1) &= \mu_v \mathbb{E}[\text{ReLU}(\|X_1\|Z)], \quad Z \sim \mathcal{N}(0, 1) \\ &= \frac{\mu_v}{2} \|X_1\| \sqrt{\frac{2}{\pi}}.\end{aligned}\tag{39}$$

Again we may take  $\mu_v = 0$  to realise  $\mu(X) = 0$ . However, in order to construct a model that is statistically not mis-specified, when using  $\sigma = \text{ReLU}$  and  $\rho = \text{u}$  it is useful to consider the case where  $\mu_v$  is non-zero (say 1). Otherwise, the model tries to describe symmetric observations  $\Gamma(X_1) = A \cdot \text{ReLU}(QX_1)$  that are equally likely negative or positive as a Gaussian distribution with a skewed non-negative mean. In order to handle non-zero  $\mu_v$ , we require evaluating additional cross-terms, given by

$$\begin{aligned}\kappa_{\sigma, \rho}(\Psi_{11}, \Psi_{22}, \Psi_{12}) &= \mathbb{E}_{(\chi, \chi')^\top \sim \mathcal{N}(\mathbf{0}, \Psi)}[\text{ReLU}(\chi) \text{u}(\chi')] \\ &= \mathbb{E}_{(\chi, \chi')^\top \sim \mathcal{N}(\mathbf{0}, \Psi)}[\chi \text{u}(\chi) \text{u}(\chi')] \\ &= \Psi_{11} \mathbb{E}[\delta(\chi) \text{u}(\chi')] + \Psi_{12} \mathbb{E}[\text{u}(\chi) \delta(\chi')] \quad (\text{multivariate Stein's lemma}) \\ &= \left( \Psi_{11} \mathbb{E}[\delta(\sqrt{\Psi_{11}} Z_1) \text{u}(\sqrt{\Psi_{22}}(Z_1 \cos \theta + Z_2 \sin \theta))] \right. \\ &\quad \left. + \Psi_{12} \mathbb{E}[\text{u}(\sqrt{\Psi_{11}}(Z_2 \cos \theta + Z_1 \sin \theta)) \delta(\sqrt{\Psi_{22}} Z_2)] \right),\end{aligned}\tag{40}$$

where  $(Z_1, Z_2)^\top \sim \mathcal{N}(0, \mathbf{I})$  and  $\cos \theta = \frac{\Psi_{12}}{\sqrt{\Psi_{22}\Psi_{11}}}$ . We have that

$$\begin{aligned}&\mathbb{E}[\delta(\sqrt{\Psi_{11}} Z_1) \text{u}(\sqrt{\Psi_{22}}(Z_1 \cos \theta + Z_2 \sin \theta))] \\ &= \frac{1}{2\pi} \int \exp\left(-\frac{1}{2}(z_1^2 + z_2^2)\right) \delta(\sqrt{\Psi_{11}} z_1) \text{u}(\sqrt{\Psi_{22}}(z_1 \cos \theta + z_2 \sin \theta)) dz_1 dz_2 \\ &= \frac{1}{2\pi\sqrt{\Psi_{11}}} \int \exp\left(-\frac{1}{2}z_2^2\right) \text{u}(\sqrt{\Psi_{22}} z_2 \sin \theta) dz_2 \\ &= \frac{1}{\sqrt{2\pi}\sqrt{\Psi_{11}}} \int \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z_2^2\right) \text{u}(z_2) dz_2 \\ &= \frac{1}{2\sqrt{2\pi}\sqrt{\Psi_{11}}}.\end{aligned}\tag{41}$$

Combining (39), (40) and (41), we have

$$\mu(X_1) \kappa_{\sigma, \rho}(\Psi_{11}, \Psi_{22}, \Psi_{12}) = \mu_v \frac{\sqrt{\Sigma_{11}}}{4\pi} \left( \sqrt{\Psi_{11}} + \frac{\Psi_{12}}{\sqrt{\Psi_{22}}} \right).\tag{42}$$

The terms involving  $\kappa_\rho$  and  $\kappa_\sigma$  are arc-cosine kernels and are given in (46) and (47).

## G Examples

### G.1 Error function example

We consider a special case where inputs are mapped to a Gaussian with a conditional expectation between  $-1$  and  $1$  through the random mapping  $\Gamma$ . We then use a Gaussian likelihood with a choice of  $R$  that maps to values between  $-1$  and  $1$ . Equivalently, we use an inverse link function that maps to values between  $-1$  and  $1$ .

Let  $p$  the pdf of a univariate standard Gaussian. Suppose input  $X$  is mapped to data  $Y = \text{erf}(WX/\sqrt{2}) + Q$  for some linear mapping  $W \in \mathbb{R}^{d \times l}$  and noise  $Q \sim \mathcal{N}(0, I)$  each with elements drawn i.i.d. from a standard Gaussian. An appropriate model is then to let  $A = \eta^2/2$  and  $R(a) = \text{erf}(a/\sqrt{2})$ . Then  $\rho(a) = 2p(a)$  and  $\sigma(a) = 2p(a) \text{erf}(a/\sqrt{2})$ .

We now invoke the general Theorem 4. In the following, we compute the individual terms in the update rule. Recall from (11), that for a particular activation function  $\zeta$ , a neural network kernel (NNK) is computed by taking the bivariate Gaussian expectation,

$$\kappa_\zeta(\Phi_{11}, \Phi_{22}, \Phi_{12}) = \mathbb{E}_{(\chi_1, \chi_2)^\top \sim \mathcal{N}(\mathbf{0}, \Phi)} [\zeta(\chi_1) \zeta(\chi_2)], \quad \Phi \triangleq \begin{pmatrix} \phi_1^\top \phi_1 & \phi_1^\top \phi_2 \\ \phi_2^\top \phi_1 & \phi_2^\top \phi_2 \end{pmatrix}.$$

It is helpful to write covariance matrices in terms of variances  $\phi_1^\top \phi_1 = \Phi_{11}$ ,  $\phi_2^\top \phi_2 = \Phi_{22}$  and covariances  $\phi_1^\top \phi_2 = \sqrt{\Phi_{11}\Phi_{22}} \cos \theta$ , where  $\theta$  is the angle between  $\phi_1$  and  $\phi_2$ . That is,

$$\Phi = \begin{pmatrix} \Phi_{11} & \sqrt{\Phi_{11}\Phi_{22}} \cos \theta \\ \sqrt{\Phi_{11}\Phi_{22}} \cos \theta & \Phi_{22} \end{pmatrix}.$$

The resulting determinant and inverse then satisfy

$$\det \Phi = \Phi_{11} \Phi_{22} \sin^2 \theta$$

$$\Phi^{-1} = \frac{1}{\Phi_{11} \Phi_{22} \sin^2 \theta} \begin{pmatrix} \Phi_{22} & -\sqrt{\Phi_{11}\Phi_{22}} \cos \theta \\ -\sqrt{\Phi_{11}\Phi_{22}} \cos \theta & \Phi_{11} \end{pmatrix}.$$

**The NNK for factor activations** A more general result is given in Tsuchida (2020, Proposition 20). For completeness, we reproduce the result here. For activation function  $\rho(a) = 2p(a)$ , we expand the 2D integral corresponding to the expectation for the NNK  $\kappa_\rho$ ,

$$\begin{aligned} \kappa_\rho(\Phi_{11}, \Phi_{22}, \Phi_{12}) &= \frac{4}{2\pi} \int \exp\left(-\frac{1}{2}(a_1^2 + a_2^2)\right) \frac{1}{2\pi\sqrt{\Phi_{11}\Phi_{22}} \sin \theta} \exp\left(-\frac{1}{2}(a_1, a_2) \Phi^{-1} (a_1, a_2)^\top\right) da_1 da_2 \\ &= \frac{2}{\pi} \int \frac{1}{2\pi\sqrt{\Phi_{11}\Phi_{22}} \sin \theta} \exp\left(-\frac{1}{2}(a_1, a_2) (\Phi^{-1} + I) (a_1, a_2)^\top\right) da_1 da_2. \end{aligned} \quad (43)$$

We now complete the square inside the argument of exp, so that we may express the integrand of (43) as a product of a bivariate Gaussian pdf and a constant.

Letting  $F^{-1} = \Phi^{-1} + I$ , we compute  $F$  as

$$\begin{aligned} F^{-1} &= \frac{1}{\Phi_{11} \Phi_{22} \sin^2 \theta} \begin{pmatrix} \Phi_{22}(1 + \Phi_{11} \sin^2 \theta) & -\sqrt{\Phi_{11}\Phi_{22}} \cos \theta \\ -\sqrt{\Phi_{11}\Phi_{22}} \cos \theta & \Phi_{11}(1 + \Phi_{22} \sin^2 \theta) \end{pmatrix} \\ \det F^{-1} &= 1 + \det \Phi^{-1} + \text{Trace} \Phi^{-1} \\ &= \frac{1 + \Phi_{11} + \Phi_{22} + \Phi_{11} \Phi_{22} \sin^2 \theta}{\Phi_{11} \Phi_{22} \sin^2 \theta} \\ F &= \frac{1}{1 + \Phi_{11} + \Phi_{22} + \Phi_{11} \Phi_{22} \sin^2 \theta} \begin{pmatrix} \Phi_{11}(1 + \Phi_{22} \sin^2 \theta) & \sqrt{\Phi_{11}\Phi_{22}} \cos \theta \\ \sqrt{\Phi_{11}\Phi_{22}} \cos \theta & \Phi_{22}(1 + \Phi_{11} \sin^2 \theta) \end{pmatrix}. \end{aligned}$$

We then rewrite (43) as

$$\begin{aligned}\kappa_\rho(\Phi_{11}, \Phi_{22}, \Phi_{12}) &= \frac{2\sqrt{\det \mathbf{F}}}{\pi\sqrt{\Phi_{11}\Phi_{22}}\sin\theta} \underbrace{\int \frac{1}{2\pi\sqrt{\det \mathbf{F}}} \exp\left(-\frac{1}{2}(a_1, a_2)\mathbf{F}^{-1}(a_1, a_2)^\top\right) da_1 da_2}_{=1} \\ &= \frac{2}{\pi\sqrt{1 + \Phi_{11} + \Phi_{22} + \Phi_{11}\Phi_{22}\sin^2\theta}} \\ &= \frac{2}{\pi\sqrt{(1 + \Phi_{11})(1 + \Phi_{22}) - \Phi_{12}^2}}\end{aligned}$$

**The NNK for chain activations** This result follows by a similar completing the square type derivation, but instead of the resulting integrand being a bivariate Gaussian density, the resulting integrand is a product of a bivariate Gaussian density with probit activations. The result then follows from Williams (1997). Concretely, the NNK  $\kappa_\sigma$  satisfies

$$\kappa_\sigma(\Phi_{11}, \Phi_{22}, \Phi_{12}) = \frac{4}{2\pi} \int \frac{\text{erf}(a_1/\sqrt{2})\text{erf}(a_2/\sqrt{2})}{2\pi\sqrt{\Phi_{11}\Phi_{22}}\sin\theta} \exp\left(-\frac{1}{2}(a_1, a_2)(\Phi^{-1} + \mathbf{I})(a_1, a_2)^\top\right) da_1 da_2.$$

Completing the square, we have

$$\begin{aligned}\kappa_\sigma(\Phi_{11}, \Phi_{22}, \Phi_{12}) &= \frac{2\sqrt{\det \mathbf{F}}}{\pi\sqrt{\Phi_{11}\Phi_{22}}\sin\theta} \underbrace{\int \frac{\text{erf}(a_1/\sqrt{2})\text{erf}(a_2/\sqrt{2})}{2\pi\sqrt{\det \mathbf{F}}} \exp\left(-\frac{1}{2}(a_1, a_2)\mathbf{F}^{-1}(a_1, a_2)^\top\right) da_1 da_2}_{=\kappa_{\text{erf}(\cdot/\sqrt{2})}(F_{11}, F_{22}, F_{12})} \\ &= \frac{2}{\pi\sqrt{(1 + \Phi_{11})(1 + \Phi_{22}) - \Phi_{12}^2}} \left( \frac{2}{\pi} \sin^{-1} \frac{F_{12}}{\sqrt{(1 + F_{11})(1 + F_{22})}} \right)\end{aligned}$$

where the last line follows from equation (11) of Williams (1997).

**The explicit kernel  $c$**  By a law of large numbers, we have that the explicit kernel is an NNK,

$$\begin{aligned}c(X_1, X_2) &= \mathbb{E}[\text{erf}(Z_1/\sqrt{2})\text{erf}(Z_2/\sqrt{2})] + 1 \quad (Z_1, Z_2) \sim \mathcal{N}\left(0, \begin{pmatrix} X_1^\top X_1 & X_1^\top X_2 \\ X_2^\top X_1 & X_2^\top X_2 \end{pmatrix}\right), \\ &= \frac{2}{\pi} \sin^{-1} \frac{X_1^\top X_2}{\sqrt{(1 + X_1^\top X_1)(1 + X_2^\top X_2)}} + 1,\end{aligned}$$

again invoking the result of Williams (1997).

**The explicit mean  $\mu$**  By a law of large numbers, the average  $\frac{1}{d} \mathbf{1}^\top \text{erf}(\mathbf{W}\mathbf{X}/\sqrt{2})$  converges to zero as  $d \rightarrow \infty$ . We therefore have that  $\mu(X) = 0$ .

## G.2 Other examples

We now investigate some other important examples. In each example, the central question is whether or not a unique fixed point exists. By Theorem 1,  $\mathbf{G}$  admits a unique fixed point if it is a contraction. It is a contraction whenever its Jacobian determinant is less than 1. The Jacobian is lower triangular, since  $0 = \frac{\partial G}{\partial \Phi_{22}} = \frac{\partial G}{\partial \Phi_{12}} = \frac{\partial G}{\partial \Phi_{11}} = \frac{\partial G}{\partial \Phi_{11}}$ , so in order to compute the Jacobian determinant, it suffices to compute the diagonal entries. These can be computed with the following identity.

**Theorem 20** (Theorem 6 of Tsuchida et al. (2021). See also Theorem 3 of Han et al. (2022).). *Suppose the absolute value of  $\zeta : \mathbb{R} \rightarrow \mathbb{R}$  is bounded by a polynomial. Let  $\dot{\zeta}$  denote the distributional (Schwartz) derivative of  $\zeta$ . Then  $\frac{\partial k_\zeta(\Phi_{11}, \Phi_{22}, \Phi_{12})}{\partial \Phi_{12}} = k_{\dot{\zeta}}(\Phi_{11}, \Phi_{22}, \Phi_{12})$  and  $\frac{\partial k_\zeta(\Phi_{11}, \Phi_{11}, \Phi_{11})}{\partial \Phi_{11}} = \mathbb{E}[(Z^2 - 1)\zeta^2(\sqrt{\Phi_{11}}Z)]/(2\Phi_{11})$ , where  $Z \sim \mathcal{N}(0, 1)$ .*

Theorem 20 allows one to compute the kernel  $k_{\sigma'}$ , where  $\sigma'$  is the derivative of  $\sigma$ , by differentiating the kernel  $k_{\sigma}$ . This is easier than computing  $k_{\sigma'}$  from scratch. There are two immediate uses for such a result. Firstly, the quantity  $k_{\sigma'}$  is needed to compute the neural tangent kernel. Secondly, and the reason the theorem is useful in our current context, is that it lets us have sufficient conditions for the update  $\mathbf{G}$  to be a contraction.

### G.2.1 Gaussian $A(\eta) = \eta^2/2$ , identity $R(a) = a$ , general $C$ , zero $\mu$

This important special case yields an  $\ell$ DEK that may be computed in closed form. Setting  $\sigma(z) = z$ , Theorem 4 and Corollary 5 say that the DEK converges to an  $\ell$ DEK with a closed-form,

$$\Psi_{ij}^{(t+1)} = \frac{1}{\lambda^2} (C_{ij} + \Psi_{ij}^{(t)}) \quad \text{and} \quad \Psi_{ij} = \frac{1}{\lambda^2} (C_{ij} + \Psi_{ij}) \implies \Psi_{ij} = \frac{C_{ij}}{\lambda^2 - 1},$$

whenever  $\lambda > 1$ , since  $\lambda > 1$  implies  $\mathbf{G}$  is a contraction. In this particular case, the  $\ell$ DEK is simply a rescaling of the kernel  $c$ .

### G.2.2 General $A$ , identity $R(a) = a$ , general $C$ , zero $\mu$

In the general setting of § 3.3, Theorem 4 and Corollary 5 yield fixed point equations for the  $\ell$ DEK that do not in general admit a closed-form,

$$\Psi_{ij}^{(t+1)} = \frac{1}{\lambda^2} \left( C_{ij} + k_{\sigma}(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}) \right) \quad \text{and} \quad \Psi_{ij} = \frac{1}{\lambda^2} \left( C_{ij} + k_{\sigma}(\Psi_{ii}, \Psi_{jj}, \Psi_{ij}) \right).$$

By Theorem 20, the  $\ell$ DEK is the fixed point of a contraction whenever  $\kappa_{\sigma}/\lambda^2 < 1$ , for which it is sufficient that  $A''$  is less than  $\lambda$ . Statistically speaking, since  $A$  acts as a cumulant generating function, this is equivalent to the largest variance of the exponential family being less than  $\lambda$ .

### G.2.3 General $A$ , general $R$ , zero $C$ , zero $\mu$

A pathological but informative example is obtained when the PSD kernel  $c$  and the cross terms  $\mu$  are chosen to be the constant zero function. In this case, from Theorem 4 and Corollary 5 we obtain

$$\Psi_{ij}^{(t+1)} = \frac{1}{\lambda^2} \kappa_{\sigma}(\Psi_{ii}^{(t)}, \Psi_{jj}^{(t)}, \Psi_{ij}^{(t)}) \quad \text{and} \quad \Psi_{ij} = \frac{1}{\lambda^2} \kappa_{\sigma}(\Psi_{ii}, \Psi_{jj}, \Psi_{ij}).$$

The  $\ell$ DEK is the fixed point of a contraction whenever  $\kappa_{\sigma}/\lambda^2 < 1$ , by Theorem 20.

Note that the (infinite  $\tau$ )  $\ell$ DEK does not depend on the input  $X_1, X_2$ , but the (finite  $\tau$ ) DEK depends on the initial guess. For a given initial guess of  $\Psi_{11}^{(1)} = \|X_1\|^2, \Psi_{22}^{(1)} = \|X_2\|^2, \Psi_{12}^{(1)} = X_1^{\top} X_2$ , solving for the  $\ell$ DEK using  $\tau$  iterations of naive fixed point iteration is exactly the same as a  $\tau$ -layer NNK (12). Therefore, the DEK is an NNK if for an arbitrary activation  $\sigma$  there exist corresponding configurations of  $A$  and  $R$ .

### G.2.4 Gaussian $A(\eta) = \eta^2/2$ , ReLU $R(a) = a \mathbf{u}(a)$ , general $C$ , zero $\mu$

Let  $\mathbf{u}$  denote the Heaviside step function, which takes values 0, 1/2 and 1 when evaluated at  $< 0$ , 0 and  $> 0$  respectively. The rectified linear unit may be written  $\text{ReLU}(a) = a \mathbf{u}(a)$ . Choosing  $A(\eta) = \eta^2/2$ ,  $\text{ReLU } R(a) = a \mathbf{u}(a)$ , we find that  $\rho$  is the Heaviside step function and  $\sigma$  is the ReLU. The corresponding kernels  $k_{\rho}$  and  $k_{\sigma}$  are known as the arc-cosine kernels of order 0 and 1, and have closed-form expressions (see Appendix I),

$$\begin{aligned} \kappa_{\mathbf{u}}(\Sigma_{11}, \Sigma_{22}, \Sigma_{12}) &= \frac{1}{2\pi} (\pi - \theta), \\ \kappa_{\text{ReLU}}(\Sigma_{11}, \Sigma_{22}, \Sigma_{12}) &= \frac{\sqrt{\Sigma_{11}\Sigma_{22}}}{2\pi} (\sin \theta + (\pi - \theta) \cos \theta), \end{aligned}$$

where  $\theta = \arccos \frac{\Sigma_{12}}{\sqrt{\Sigma_{11}\Sigma_{22}}}$ .

Fixed points for  $\Psi_{11}$  and  $\Psi_{22}$  can be computed in closed-form provided  $\lambda^2 < 1/2$ ,

$$\Psi_{ii}^{(t+1)} = \frac{1}{2\lambda^2} (C_{ii} + \Psi_{ii}^{(t)}) \quad \text{and} \quad \Psi_{ii} = \frac{1}{2\lambda^2} (C_{ii} + \Psi_{ii}) \implies \Psi_{ii} = \frac{C_{ii}}{2\lambda^2 - 1}.$$

However,  $\Psi_{12}$  cannot be computed in closed-form. We leave the analysis for determining whether G results in a contraction for future work. Nevertheless, we may still compute using the result in Theorem 4 without violating any assumptions.

Interestingly, in this setting,  $\rho = \dot{\sigma}$  almost everywhere and the DEK iterates very closely resemble NTK iterates. There are three differences in calculating the DEK and the NTK. Firstly, the DEK uses  $c$  where the NTK uses  $\Theta^{(t)}$ . Secondly, the DEK uses  $\Psi^{(t)}$  as an input to  $\Sigma^{(t+1)}$  and  $\dot{\Sigma}^{(t+1)}$ , whereas the NTK uses  $\Sigma^{(t)}$ . Finally, the DEK may be initialised at any guess  $\Psi^{(1)}$ , whereas the NTK must be initialised at  $X_1^\top X_2$ .

### G.2.5 Gaussian $A(\eta) = \eta^2/2$ , ReLU $R(a) = a \mathbf{u}(a)$ , arc-cosine $c$ and corresponding $\mu$

We now describe a setting that we found practically useful in our experiments (see § 4.2.). We use the setting described in § G.2.4, but without the assumption that  $\mu(X) = 0$ . For the features  $\Gamma(X)$  of the kernel  $c(X_1, X_2)$ , we choose  $\Gamma(X) = \mu_v \text{ReLU}(QX)$ , where  $\mu_v \in \mathbb{R}$  is a hyperparameter and  $Q$  is a  $d \times l$  matrix with entries drawn independently from the standard Gaussian distribution, resulting in an arc-cosine kernel for  $c$ . The mean function  $\mu$  and the cross terms  $\kappa_{\sigma, \rho}$  admit closed-form expressions, as given in Appendix F.1. The resulting DEK can represent deep arc-cosine kernels when  $\mu_v$  is zero, and resembles (but is not the same as) an NTK with extra cross-terms otherwise.

## H Other considerations

### H.1 Why the expected negative log posterior?

We may frame our optimisation objective in terms of exponential family PCA (Collins et al., 2001; Mohamed et al., 2008). Given a dataset  $\{Y_s\}_{s=1}^N$  of  $N$  examples, exponential family PCA models observation  $Y_s \in \mathbb{Y}^d$  as following a factored exponential family with canonical parameter  $V\phi_s$ , for some basis  $V$  and latent  $\phi_s \in \mathbb{R}^m$ . The resulting graphical model is shown in Figure 3a. A maximum a posteriori estimate is

$$\phi_s^* \triangleq \arg \min -\log p(\phi_s | Y_s) = \arg \min -\log \int p(Y_s | V, \phi_s) p(V) p(\phi_s) dV, \quad (44)$$

in which the basis  $V$  is marginalised before the evaluation of the logarithm.

Our objective (17) differs from (44) in two respects. Firstly, we generalise the canonical parameter  $V\phi_s$  so that a nonlinearly parameterised canonical parameter  $R(V\phi_s)$  is used. Secondly and more critically, the order of the logarithm and the expectation is swapped. This may be understood by examining a variational lower bound (VLB) of the posterior. Note that the VLB has been used for MAP estimation in similar contexts (Kingma & Welling, 2014), and can be seen as a regularised or penalised variant of the ELBO. Let  $\Phi = (\phi_1, \dots, \phi_N)$ . For any density  $q(V)$  which ostensibly approximates  $p(V | Y, \Phi)$ , the log model evidence decomposes into a sum of a KL divergence and an ELBO,

$$\begin{aligned} \log p(Y | \Phi) &= \text{KL}(q(V) \| p(V | Y, \Phi)) + \mathbb{E}_{V \sim q} \log \frac{p(V, Y | \Phi)}{q(V)}, \quad \text{so that} \\ -\mathbb{E}_{V \sim q} \log p(\Phi | V, Y) &= \text{KL}(q(V) \| p(V | Y, \Phi)) - \log p(\Phi | Y) + \mathbb{E}_{V \sim q} \log p(V, Y) - \mathbb{E}_{V \sim q} \log q(V). \end{aligned}$$

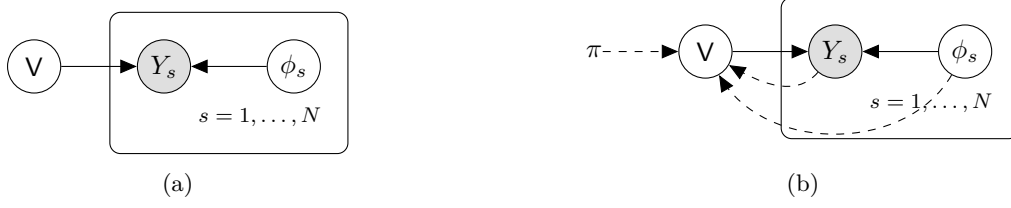
By minimising the left hand side with respect to  $\Phi$ , we are maximising the log model evidence minus the KL divergence. By selecting hyperparameters  $\pi$  of the variational density  $q$  over  $V$ , we alter our approximate posterior.

Note a clashing nomenclature between EM-algorithm and variational inference — where the marginalised variable  $V$  is called a latent variable — against an unsupervised dimensionality reduction setting — where the low dimensional representation  $\phi_s$  is called a latent variable.

### H.2 Scaling and parameterisation of weight distributions

It is widely appreciated that the prior over  $W$  in the Bayesian setting (MacKay, 1998, §11.1) and the initialisation of  $W$  in the gradient-flow setting play an role in directing the limiting behaviour of the neural network (Sohl-Dickstein et al., 2020). On the one hand, convenient parameterisations and choices of prior and initial distributions lead to tractable large width limits. On the other hand, while limiting models can outperform their finite width counterparts in small data regimes (Arora et al., 2020), GPs in general are most often outperformed by deep learning models for many problems of interest. This might suggest that the tractable limits are the “wrong” ones to analyse if one seeks to explain the success stories of deep learning (Chizat et al., 2019; Woodworth et al., 2020). Other works consider more general heavy-tailed (Der & Lee, 2005; Peluchetti et al., 2020; Favaro et al., 2021; 2022) or differently scaled priors, but it is not yet clear whether these models can more accurately emulate deep learning models.

In our work in particular, the scaling of the prior with precision  $\sqrt{md}$  (less than the  $m$  that might often be expected, since  $d < m$ ) in (17) was crucial for finding a tractable limit. Independently of whether this limiting regime represents any meaningful feature representation, our analysis is valuable because (1) DEKs are better than or competitive with other neural network kernel models in the settings that we tried, (2) we are the first to place deep neural network related kernels in a more fundamental footing of statistical estimation and optimisation, and (3) our analysis describes a limiting invariant of SGD.



### H.3 Implicit differentiation

From Corollary 5, we have that the  $\ell$ DEK  $\Psi$  satisfies  $\Psi = G(\Psi)$ . Suppose  $\Psi$  depends on  $v$ -dimensional hyperparameter  $\zeta \in \mathbb{R}^v$ , such as the weight and bias variance (see Footnotes 1 and 2), or a hyperparameter of  $R$ . If  $G$  is continuously differentiable, the implicit function theorem says

$$\underbrace{\frac{d\Psi}{d\zeta}}_{3 \times v} = \underbrace{\frac{\partial G(\Psi)}{\partial \zeta}}_{3 \times v} + \underbrace{\frac{\partial G(\Psi)}{\partial \Psi}}_{3 \times 3} \underbrace{\frac{d\Psi}{d\zeta}}_{3 \times v} \implies \left(1 - \frac{\partial G(\Psi)}{\partial \Psi}\right) \frac{d\Psi}{d\zeta} = \frac{\partial G(\Psi)}{\partial \zeta},$$

which may be solved for  $\frac{d\Psi}{d\zeta}$  using a backslash operator. This derivative may be used for gradient-based hyperparameter selection. For example, if the  $\ell$ DEK were to be used as the covariance function of a Gaussian process, one could perform type II maximum marginal likelihood to compute point estimates for  $\zeta$ . This implicit differentiation mirrors the finite-width counterpart, the DEQ (Bai et al., 2019). We leave its empirical investigation for future work.



## I Arc-cosine kernels via derivatives

While the Dirac distribution is not a function and therefore cannot be used as an activation function in finite-width networks, it does arise as the derivative of NNKs with Heaviside activations, by Theorem 20. With an abuse of notation that extends the usual operation of integrating against a Dirac delta distribution, we may understand an expectation involving Dirac delta distributions as a limiting expectation involving nascent delta functions. We may evaluate the corresponding NNK as follows.

$$\begin{aligned}
& \kappa_\delta(\Sigma_{11}, \Sigma_{22}, \Sigma_{12}) \\
&= \mathbb{E}[\delta(\chi)\delta(\chi')] \\
&= \mathbb{E}[\delta(\sqrt{\Sigma_{11}}Z_1)\delta(\sqrt{\Sigma_{22}}(Z_1\rho + Z_2\sqrt{1-\rho^2}))], \quad (Z_1, Z_2)^\top \sim \mathcal{N}(0, I), \quad \rho = \Sigma_{12}/\sqrt{\Sigma_{11}\Sigma_{22}} \\
&= \frac{1}{\sqrt{\Sigma_{11}}} \frac{1}{\sqrt{2\pi}} \int \delta(\sqrt{\Sigma_{22}}z_2\sqrt{1-\rho^2})p(z_2)dz_2, \quad p \text{ is pdf of standard Gaussian} \\
&= \frac{1}{\sqrt{\Sigma_{11}\Sigma_{22}(1-\rho^2)}} \frac{1}{2\pi} \\
&= \frac{1}{2\pi\sqrt{\Sigma_{11}\Sigma_{22} - \Sigma_{12}^2}}.
\end{aligned} \tag{45}$$

Note the singularity whenever the Gaussian distribution is degenerate, i.e.  $\Sigma_{11} = \Sigma_{22} = \Sigma_{12}$ , which is an instance of the more general undefinedness of a product of Dirac delta distributions.

The NNK corresponding with Heaviside activations  $u$  was first evaluated using a geometric argument by Sheppard (1899), and is given by

$$\begin{aligned}
& \kappa_u(\Sigma_{11}, \Sigma_{22}, \Sigma_{12}) \\
&= \mathbb{E}[u(\chi)u(\chi')] \\
&= \frac{1}{2\pi}(\pi - \theta), \quad \theta = \arccos \frac{\Sigma_{12}}{\sqrt{\Sigma_{11}\Sigma_{22}}}.
\end{aligned} \tag{46}$$

The NNK (46) was generalised to activations of the form  $u(z)z^p$  for positive integers  $p$  by Cho & Saul (2009). Of particular relevance is the case  $p = 1$ , in which case the activation function is ReLU and

$$\begin{aligned}
& \kappa_{\text{ReLU}}(\Sigma_{11}, \Sigma_{22}, \Sigma_{12}) \\
&= \mathbb{E}[\text{ReLU}(\chi)\text{ReLU}(\chi')] \\
&= \frac{\sqrt{\Sigma_{11}\Sigma_{22}}}{2\pi}(\sin \theta + (\pi - \theta)\cos \theta).
\end{aligned} \tag{47}$$

Note that (45)–(47) represent a sequence of derivatives, since the Dirac delta distribution, Heaviside function and ReLU represent a sequence of distributional derivatives. More concretely, by Theorem 20,

$$\frac{\partial^2 \kappa_{\text{ReLU}}}{\partial \Sigma_{12}^2} = \frac{\partial \kappa_u}{\partial \Sigma_{12}} = \kappa_\delta,$$

as can be otherwise verified.

## J Experiments

### J.1 Measuring finite-width effects

We consider elements of an input space  $\mathbb{X}$  which are 100 evenly spaced points over  $[-5, 5]^2$ . This results in an input matrix of size  $100 \times 2$ . We compute two  $100 \times 100$  kernel matrices with  $ij$ th element:  $\Psi_{ij}$  (calculated to high tolerance using a fixed point solver) and  $k_d^{(t)}(X_i, X_j)$  (calculated using SGD). Finite features  $\Gamma$  are chosen to be  $\Gamma = \mathbf{T}X$ , where  $\mathbf{T} \in \mathbb{R}^{d \times m}$  is a zero-mean Gaussian random matrix. This results in a linear kernel  $c(X_1, X_2) = X_1^\top X_2$ . We set  $t = 400$ ,  $\lambda = 6$  and use a step length of  $\alpha^{(t)} = \frac{1}{\lambda} \sqrt{d/m}$ . We vary  $d$  between 5 and 500 in steps of 5 and choose  $m = d^{3/2}$ . We also provide the CKA between the (finite- $d$ , finite- $\tau$ ) ffDEK and a squared exponential kernel  $A \exp(-\|X_1 - X_2\|_2^2/2)$  for control, where the scaling parameter  $A$  is the largest value in the (infinite- $d$ , infinite- $\tau$ )  $\ell$ DEK matrix.

### J.2 Inference using the DEK

The hyperparameter grid over which `GridSearchCV` operates is given in table 4.

Hyperparameter	Present in	Values
Data scale (see footnote 1)	NTK, NNK, DEK, SEK	$\{0.5, 1, 2, 4\}$
KRR regularisation strength	NTK, NNK, DEK, SEK	$\{0.05, 0.1, 0.5\}$
Input augmented bias (see footnote 2)	NTK, NNK, DEK	$\{-1.0, -0.1, 0.0, 0.1, 1.0\}$
Number of iterations / layers $T$	NTK, NNK	$\{2, 3, 4, 5\}$
Number of iterations / layers $T$	DEK	$\{2, 3, 4, 5, \infty\}$
Inner regularisation strength $\lambda$	DEK	$\{1, 2, 4\}$
Cross-term strength $\mu_v$	DEK	$\{0, 0.1, 0.5, 1, 2\}$
Lengthscale	SEK	$\{0.5, 1, 2, 4, 8, 16\}$

Table 4: Search space for `GridSearchCV`. We use Anderson acceleration to compute the DEK when  $T = \infty$  and there are less than 500 points in the dataset.