

Notes on Neural Tangent Kernels

Luigi Del Debbio and Amedeo Chiefa

Higgs Centre for Theoretical Physics, School of Physics and Astronomy,
Peter Guthrie Tait Road, Edinburgh EH9 3 FD, United Kingdom.

May 8, 2025

Abstract

Summary of results for NTKs. These notes follow the notation introduced in Ref. [1].

1 Definitions

1.1 Neural Networks

We consider the case of neural networks (NNs) with $L + 1$ layers, labelled by the index $\ell = 0, \dots, L$. The number of neurons in layer ℓ is denoted by n_ℓ . The $\ell = 0$ and $\ell = L$ layers are the input and output layers respectively, so that a NN parametrizes an element of

$$\mathcal{F} = \{f : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_L}\} . \quad (1)$$

The space of inputs in this case is $\mathcal{X} = \mathbb{R}^{n_0}$, while the space of outputs is $\mathcal{Y} = \mathbb{R}^{n_L}$. The number of parameters used in the parametrization is the total number of weights and biases,

$$P = \sum_{\ell=0}^{L-1} (n_\ell + 1)n_{\ell+1} . \quad (2)$$

The weights connecting the layer ℓ to the layer $\ell + 1$ are denoted $W^{(\ell+1)}$ and the biases in layer $\ell + 1$ are denoted $b^{(\ell+1)}$. With this convention

$$W^{(\ell+1)} \in \mathbb{R}^{n_{\ell+1} \times n_\ell} , \quad b^{(\ell+1)} \in \mathbb{R}^{\ell+1} . \quad (3)$$

When referring to the whole set of parameters, we use the symbol

$$\theta = \left\{ W^{(\ell)}, b^{(\ell)}; \ell = 0, \dots, L - 1 \right\} , \quad (4)$$

while individual parameters are denoted θ_μ , with $\mu = 1, \dots, P$.

The input space is $\mathcal{X} = \mathbb{R}^{n_0}$ and we consider a probability distribution p^{in} over \mathcal{X} , which induces a scalar product in \mathcal{F} ,

$$\langle f, g \rangle_{\text{in}} = \mathbb{E}_{x \sim p^{\text{in}}} [f(x)^T g(x)] \quad (5)$$

$$= \int dx p^{\text{in}}(x) f_i(x) g_i(x). \quad (6)$$

For a discrete set of datapoints \mathcal{D} ,

$$p^{\text{in}}(x) = \frac{1}{|\mathcal{D}|} \sum_{\alpha \in \mathcal{D}} \delta(x - x_\alpha), \quad (7)$$

where the index α labels the elements of the dataset. For a given set of parameters, the function parametrized by the NN, which we denote f_θ , is given by the preactivation function of the output layer L ,

$$f(x; \theta) = \phi^{(L)}(x). \quad (8)$$

The preactivations $\phi^{(\ell)}$ and the activations $\rho^{(\ell)}$ are defined recursively as

$$\rho_i^{(0)}(x) = x_i, \quad (9)$$

$$\phi_i^{(\ell+1)}(x) = \sum_{j=1}^{n_\ell} W_{ij}^{(\ell+1)} \rho_j^{(\ell)}(x) + b_i^{(\ell+1)}, \quad (10)$$

$$\rho_i^{(\ell)}(x) = \rho(\phi_i^{(\ell)}(x)), \quad (11)$$

where ρ is the activation function. The function

$$F^{(L)} : \mathbb{R}^P \rightarrow \mathcal{F} \quad (12)$$

$$\theta \mapsto F^{(L)}(\theta) = f_\theta \quad (13)$$

is called the NN realization function. We also introduce the (dual) space of linear functionals acting on \mathcal{F} ,

$$\mathcal{F}^* = \{\mu : \mathcal{F} \rightarrow \mathbb{R}; \mu = \langle d, \cdot \rangle_{\text{in}}, d \in \mathcal{F}\}, \quad (14)$$

and a loss function

$$\mathcal{L} : \mathcal{F} \times \mathcal{Z} \rightarrow \mathbb{R}. \quad (15)$$

Here \mathcal{Z} denotes the space of datapoints, which we can identify for now with $\mathcal{X} \times \mathcal{Y}$.

1.2 Kernel Gradient

A kernel K is defined as a symmetric function

$$K : \mathbb{R}^{n_0} \times \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_L \times n_L}, \quad (16)$$

such that $K(x, x') = K(x', x)^T$. Using the kernel, a bilinear map is defined on \mathcal{F}

$$\langle f, g \rangle_K = \mathbb{E}_{x, x' \sim p^{\text{in}}} [f(x)^T K(x, x') g(x')] . \quad (17)$$

The kernel is positive definite if

$$\|f\|_{p^{\text{in}}} > 0 \implies \|f\|_K > 0. \quad (18)$$

Given a kernel K , we define a map

$$\Phi_K : \mathcal{F}^* \rightarrow \mathcal{F} \quad (19)$$

$$\mu \mapsto f_\mu, \quad (20)$$

such that

$$f_{\mu, i}(x) = \langle d, K_i(x, \cdot) \rangle_{p^{\text{in}}} \quad (21)$$

$$= \sum_{j=1}^{n_L} \int dx' p^{\text{in}}(x') d_j(x') K_{ij}(x, x'), \quad (22)$$

where $d \in \mathcal{F}$, and $\mu = \langle d, \cdot \rangle_{p^{\text{in}}}$.

The cost function for a given dataset depends on the values of the function f at the datapoints x_α , so that we can write

$$C_{\mathcal{D}}(f) = \widehat{C}(\mathbf{f}), \quad (23)$$

where \mathbf{f} is the finite-dimensional vector made of the values $f_\alpha = f(x_\alpha)$ for $\alpha = 1, \dots, |\mathcal{D}|$. The variation of the cost function at a point $f_0 \in \mathcal{F}$ is

$$\delta C_{\mathcal{D}}|_{f_0} = C_{\mathcal{D}}(f_0 + \delta f) - C_{\mathcal{D}}(f_0) \quad (24)$$

$$= \sum_{\alpha \in \mathcal{D}} \left. \frac{\partial \widehat{C}}{\partial f_\alpha} \right|_{f_0} \delta f_\alpha \quad (25)$$

$$= \partial_f^{\text{in}} C|_{f_0} \delta f \quad (26)$$

$$= \langle d|_{f_0}, \delta f \rangle_{p^{\text{in}}}, \quad (27)$$

where $\delta f \in \mathcal{F}$, $\delta f_\alpha = \delta f(x_\alpha)$ and $\partial_f^{\text{in}} C|_{f_0} \in \mathcal{F}^*$.

Using the kernel K , the *kernel gradient* $\nabla_K C|_{f_0}$ is defined as $\Phi_K \left(\partial_f^{\text{in}} C|_{f_0} \right)$, *i.e.*

$$\nabla_K C|_{f_0}(x) = \frac{1}{|\mathcal{D}|} \sum_{\alpha \in \mathcal{D}} K(x, x_\alpha) \left. \frac{\partial \hat{C}}{\partial f_\alpha} \right|_{f_0} . \quad (28)$$

A function $f(t)$ follows the kernel gradient descent if

$$\partial_t f(t) = - \nabla_K C|_{f(t)} . \quad (29)$$

During the kernel descent, the cost function variation is given by

$$\partial_t C|_{f(t)} = - \langle d|_{f(t)}, \nabla_K C|_{f(t)} \rangle_{p^{\text{in}}} = - \| d|_{f(t)} \|_K^2 . \quad (30)$$

More generally, for a function \hat{O} that depends on the values of f on the points in the dataset,

$$\partial_t \hat{O}|_{f(t)} = - \left\langle \frac{\partial \hat{O}}{\partial f}, \frac{\partial \hat{C}}{\partial f} \right\rangle_K . \quad (31)$$

2 Gradient Descent for NNPfD

We work in a basis where the data covariance C_Y is diagonal. Such a choice has two main advantages:

1. The data points in this basis are statistically independent, so that when the data is split, e.g. between training and validation sets, the points in each subset are not correlated.
2. Working in a basis where the data covariance is diagonal, requires to compute the eigenvectors and eigenvalues of C_Y , thereby automatically checking for a potentially large condition number.

Note that, in case we need to apply kinematic cuts, these can be applied before rotating to the diagonal basis. The diagonal basis is then defined by determining the eigensystem of the reduced covariance matrix, i.e. the matrix obtained by removing the lines and columns that correspond to data points that do not pass the kinematic cuts.

The loss function is defined to be the Mean Square Error (MSE),

$$\mathcal{L} = \frac{1}{2} (y - T)^T C_Y^{-1} (y - T) , \quad (32)$$

$$= \frac{1}{2} \| y - T \|_{C_Y}^2 , \quad (33)$$

where

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_I \\ \vdots \\ y_{N_{\text{dat}}} \end{pmatrix}, \quad \text{and} \quad T = \begin{pmatrix} T_1 \\ \vdots \\ T_I \\ \vdots \\ T_{N_{\text{dat}}} \end{pmatrix} \quad (34)$$

are vectors containing the data points y_I and their respective theoretical predictions T_I , where $I = 1, \dots, N_{\text{dat}}$. The loss function can be written using explicit indices,

$$\mathcal{L} = \frac{1}{2} \sum_I \frac{1}{\sigma_I^2} (y_I - T_I)^2. \quad (35)$$

However, here and in what follows, we try to minimize the usage of explicit indices.

Denoting by t the training time, and by θ_μ , with $\mu = 1, \dots, P$, the parameters used to specify the functional form of f , gradient descent is defined by the flow equation,

$$\frac{d}{dt} \theta_{\mu,t} = - \sum_{\nu} \lambda_{\mu\nu} \nabla_{\nu} \mathcal{L}_t, \quad (36)$$

where we added a suffix \mathcal{L}_t to emphasise that the loss function is evaluated using the function f_t at training time t . We introduced the tensor $\lambda_{\mu\nu}$ to take into account the case of a non-trivial metric in the space of parameters. We will use $\lambda_{\mu\nu} = \lambda_{\mu} \delta_{\mu\nu}$, unless explicitly stated. The gradient on the RHS of Eq. 36 above is

$$\nabla_{\mu} \mathcal{L}_t = - (\nabla_{\mu} f_t)^T \left(\frac{\partial T}{\partial f} \right)_t^T C_Y^{-1} \epsilon_t, \quad (37)$$

where we introduced $\epsilon_t = y - T_t$, to denote the difference between the data and the theoretical prediction obtained from f_t . The quantities that depend on the training time have an explicit suffix t . Using the chain rule yields

$$\frac{d}{dt} f_t = \sum_{\mu} (\nabla_{\mu} f_t) \frac{d}{dt} \theta_{\mu,t} \quad (38)$$

$$= \left[\sum_{\mu,\nu} \lambda_{\mu\nu} (\nabla_{\mu} f_t) (\nabla_{\nu} f_t)^T \right] \left(\frac{\partial T}{\partial f} \right)_t^T C_Y^{-1} \epsilon_t. \quad (39)$$

Introducing the Neural Tangent Kernel (NTK),

$$\Theta_t = \sum_{\mu,\nu} \lambda_{\mu\nu} (\nabla_{\mu} f_t) (\nabla_{\nu} f_t)^T, \quad (40)$$

or equivalently using explicit indices for once,

$$\Theta_{\delta_1 \delta_2, t} = \sum_{\mu\nu} \lambda_{\mu\nu} \nabla_\mu f_{\delta_1, t} \nabla_\nu f_{\delta_2, t}. \quad (41)$$

Note that the NTK is a tensor in the space of the grid points on which the function f is evaluated, *i.e.* $f_\delta = f(x_\delta)$.

The flow equation becomes

$$\frac{d}{dt} f_t = \Theta_t \left(\frac{\partial T}{\partial f} \right)_t^T C_Y^{-1} \epsilon_t. \quad (42)$$

This flow equation is completely general, for any parametrization and any dependence of the theoretical predictions on the function f .

Let us now focus on the case where the theoretical prediction is linear in f ,

$$T_I = (\text{FK})_{I\alpha} f_\alpha, \quad (43)$$

so that

$$\left(\frac{\partial T_I}{\partial f_\alpha} \right) = (\text{FK})_{I\alpha}, \quad (44)$$

is actually independent of t . Note in passing that (FK) is not a square matrix, and hence not invertible.

The loss function is a quadratic form in f ,

$$2\mathcal{L} = f^T M f - f^T (\text{FK})^T C_Y^{-1} y - y^T C_Y^{-1} (\text{FK}) f + y^T C_Y^{-1} y. \quad (45)$$

We can restrict f to the subspace orthogonal to $\ker(\text{FK})$; the components in $\ker(\text{FK})$ do not enter the loss function and therefore are unconstrained. In this subspace, the matrix

$$M = (\text{FK})^T C_Y^{-1} (\text{FK}) \quad (46)$$

is symmetric and positive definite, and hence invertible. Rewriting

$$f = \chi + M^{-1} (\text{FK})^T C_Y^{-1} y, \quad (47)$$

yields

$$2\mathcal{L} = \chi^T M \chi + y^T C_Y^{-1} [1 - (\text{FK}) M^{-1} (\text{FK})^T C_Y^{-1}] y. \quad (48)$$

Therefore the absolute minimum is obtained for

$$f^* = M^{-1} (\text{FK})^T C_Y^{-1} y \implies \mathcal{L}^* = \frac{1}{2} y^T C_Y^{-1} [1 - (\text{FK}) M^{-1} (\text{FK})^T C_Y^{-1}] y. \quad (49)$$

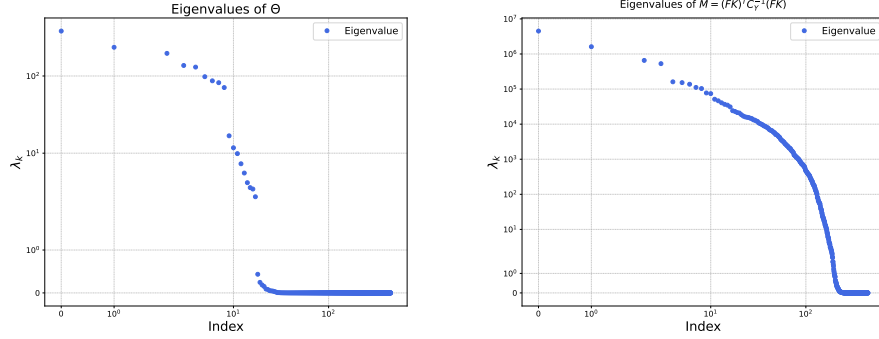


Figure 1: Eigenvalues of the NTK (left) and of the matrix M (right) in logarithmic scale.

Note that the minimum of the loss function is *not* zero. However, in a level-0 closure test, $y = (\text{FK})f^{\text{in}}$, and then $f^* = f^{\text{in}}$ and $\mathcal{L}^* = 0$. The absolute minimum of the loss yields the correct input function and a vanishing loss function. This is something we always stress in NNPfD. If we push the training for long enough in L0, we can get a vanishing loss. The flow equation for linear data becomes

$$\frac{d}{dt}f_t = \Theta_t (\text{FK})^T C_Y^{-1} \epsilon_t. \quad (50)$$

Evolution of f_t . In order to solve for the evolution of f_t , we start by finding the eigenvectors and eigenvalues of Θ ,

$$\Theta v_i = \lambda_i v_i. \quad (51)$$

The eigenvalues of the NTK, together with those of the matrix M , are displayed in Fig. 1. We rewrite the flow equation, Eq. 50, as

$$\frac{d}{dt}f_t = b - \Theta M f_t, \quad (52)$$

where

$$b = \Theta(\text{FK})^T C_Y^{-1} y, \quad (53)$$

while the matrix M has been introduced above. Using the eigendecomposition $M = R D R^T$, where R is an orthogonal matrix, we introduce

$$\tilde{f}_t = D^{1/2} R^T f_t, \quad \tilde{b} = D^{1/2} R^T b, \quad (54)$$

so that

$$\frac{d}{dt}\tilde{f}_t = \tilde{b} - D^{1/2}R^T\Theta RD^{1/2}\tilde{f}_t \quad (55)$$

$$= \tilde{b} - \tilde{H}\tilde{f}_t. \quad (56)$$

The eigenvalues of the matrix \tilde{H} are shown in Fig. 2.

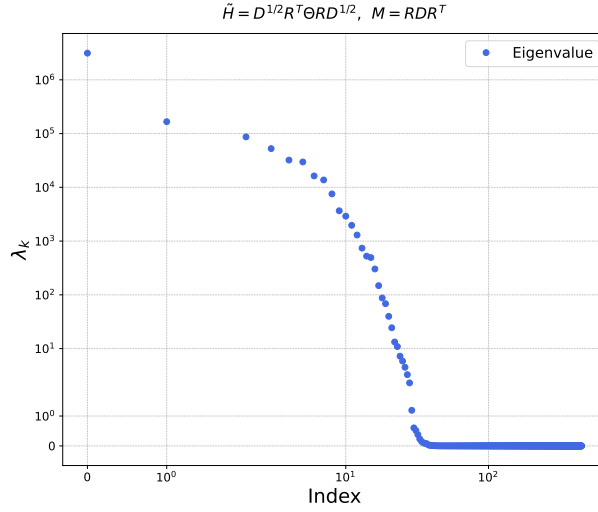


Figure 2: Eigenvalues of \tilde{H} in logarithmic scale.

The solution, assuming that \tilde{H} is independent of the flow time, is

$$\tilde{f}_t = e^{-\tilde{H}t}\tilde{f}_0 + (1 - e^{-\tilde{H}t})\tilde{H}^{-1}\tilde{b}. \quad (57)$$

Taking the limit $t \rightarrow \infty$, we find

$$f_\infty = M^{-1}(\text{FK})^T C_Y^{-1}y = f^*, \quad (58)$$

and thus

$$\epsilon_\infty = [1 - (\text{FK})M^{-1}(\text{FK})^T C_Y^{-1}]y. \quad (59)$$

A few lines of algebra yield

$$\frac{1}{2}\|\epsilon_\infty\|_{C_Y}^2 = \mathcal{L}^*, \quad (60)$$

in agreement with Eq. 49. As before, if we consider a level-0 test,

$$f_\infty = M^{-1}(\text{FK})^T C_Y^{-1}(\text{FK})f^{\text{in}} = f^{\text{in}}. \quad (61)$$

Evolution for ϵ_t . We can rewrite Eq. 50 as

$$\begin{aligned}\frac{d}{dt}\epsilon_t &= -(\text{FK})\frac{d}{dt}f_t \\ &= -(\text{FK})\Theta(\text{FK})^T C_Y^{-1}\epsilon_t \\ &= -(\text{FK})\Theta(\text{FK})^T R_Y D_Y R_Y^T \epsilon_t,\end{aligned}\tag{62}$$

where we decompose $C_Y^{-1} = R_Y D_Y R_Y^T$. Following the derivation for f_t , we define

$$\tilde{\epsilon}_t = D_Y^{1/2} R_Y^T \epsilon_t,\tag{63}$$

so that

$$\frac{d}{dt}\tilde{\epsilon}_t = -\tilde{H}_\epsilon \tilde{\epsilon}_t,\tag{64}$$

where

$$\tilde{H}_\epsilon = D_Y^{1/2} R_Y^T (\text{FK}) \Theta (\text{FK})^T R_Y D_Y^{1/2}\tag{65}$$

is a symmetric, positive definite, matrix, and therefore

$$\epsilon_t = R_Y D_Y^{-1/2} e^{-\tilde{H}_\epsilon t} D_Y^{1/2} R_Y^T \epsilon_0.\tag{66}$$

Note that

$$\begin{aligned}\tilde{H}_\epsilon \tilde{\epsilon}_\infty &= D_Y^{1/2} R_Y^T (\text{FK}) \Theta \underbrace{(\text{FK})^T C_Y^{-1} [1 - (\text{FK}) M^{-1} (\text{FK})^T C_Y^{-1}]}_{(*)} y, \\ (*) &= (\text{FK})^T C_Y^{-1} - \underbrace{(\text{FK})^T C_Y^{-1} (\text{FK}) M^{-1} (\text{FK})^T C_Y^{-1}}_{=1!} = 0,\end{aligned}\tag{67}$$

and therefore the evolution operator for $\tilde{\epsilon}$, \tilde{H}_ϵ , has a zero mode, which is not driven to zero by the gradient flow. The zero mode is exactly ϵ_∞ , *i.e.* the residual error at the end of the gradient flow. All other components of ϵ_t vanish exponentially with t . This result is consistent with the result above for the evolution of f_t . Now we need to check them numerically!

Consistency of the two evolutions. As a cross check of our calculation, we verify that the evolution of ϵ_t is consistent with the evolution of f_t . Starting from Eq. 63,

$$\tilde{\epsilon}_t = D_Y^{1/2} R_Y^T \left(y - (\text{FK}) R D^{-1/2} \tilde{f}_t \right).\tag{68}$$

We compute separately the time-independent and the time-dependent contributions to $\tilde{\epsilon}_t$ in the RHS of Eq. 68. Starting from the time-independent terms, we have

$$\begin{aligned}
y - (\text{FK})RD^{-1/2}\tilde{H}^{-1}\tilde{b} &= \\
&= y - (\text{FK}) \underbrace{RD^{-1/2} \left[D^{-1/2} R^T \Theta^{-1} \right]}_{=M^{-1}} \underbrace{RD^{-1/2}}_{=1} \left[D^{1/2} R^T \Theta (\text{FK})^T C_Y^{-1} y \right] \\
&= \epsilon_\infty,
\end{aligned} \tag{69}$$

and therefore

$$D_Y^{1/2} R_Y^T \left(y - (\text{FK})RD^{-1/2}\tilde{H}^{-1}\tilde{b} \right) = \tilde{\epsilon}_\infty. \tag{70}$$

The time-independent component of $\tilde{\epsilon}_t$ yields $\tilde{\epsilon}_\infty$, as expected. The time-dependent terms in $(\text{FK})RD^{-1/2}\tilde{f}_t$ are

$$(\text{FK})RD^{-1/2}e^{-\tilde{H}t} \left(\tilde{H}^{-1}\tilde{b} - \tilde{f}_0 \right). \tag{71}$$

Focussing on the first term above,

$$\begin{aligned}
(\text{FK})RD^{-1/2}e^{-\tilde{H}t}\tilde{H}^{-1}\tilde{b} &= \\
&= (\text{FK})RD^{-1/2}\tilde{H}^{-1}\tilde{b} + \sum_{k>0} (\text{FK})RD^{-1/2} \frac{(-t)^k}{k!} \left(D^{1/2} R^T \Theta RD^{1/2} \right)^{k-1} D^{1/2} R^T b \\
&= (\text{FK})RD^{-1/2}\tilde{H}^{-1}\tilde{b} + \sum_{k>0} \frac{(-t)^k}{k!} \underbrace{(\text{FK}) (\Theta M)^{k-1}}_{(*)} b.
\end{aligned} \tag{72}$$

Note that

$$(*) = (\text{FK}) \underbrace{\left[\Theta (\text{FK})^T C_Y^{-1} (\text{FK}) \right] \left[\Theta (\text{FK})^T C_Y^{-1} (\text{FK}) \right] \dots}_{k-1 \text{ factors}} \tag{73}$$

$$= \underbrace{\left[(\text{FK}) \Theta (\text{FK})^T C_Y^{-1} \right] \left[(\text{FK}) \Theta (\text{FK})^T C_Y^{-1} \right] \dots (\text{FK})}_{k-1 \text{ factors}} \tag{74}$$

$$= \left[(\text{FK}) \Theta (\text{FK})^T C_Y^{-1} \right]^{k-1} (\text{FK}), \tag{75}$$

and therefore

$$(\text{FK})RD^{-1/2}e^{-\tilde{H}t}\tilde{H}^{-1}\tilde{b} = \sum_k \frac{(-t)^k}{k!} \left[(\text{FK}) \Theta (\text{FK})^T C_Y^{-1} \right]^k y - \epsilon_\infty. \tag{76}$$

Hence the contribution to $\tilde{\epsilon}_t$ from this term

$$\begin{aligned}
& \sum_k \frac{(-t)^k}{k!} D_Y^{1/2} R_Y^T \left[(\text{FK}) \Theta (\text{FK})^T R_Y D_Y^{1/2} D_Y^{1/2} R_Y^T \right]^k y = \\
& = \sum_k \frac{(-t)^k}{k!} \left[D_Y^{1/2} R_Y^T (\text{FK}) \Theta (\text{FK})^T R_Y D_Y^{1/2} \right]^k D_Y^{1/2} R_Y^T y \\
& = e^{-\tilde{H}_\epsilon t} (\tilde{y} - \tilde{\epsilon}_\infty) .
\end{aligned} \tag{77}$$

The second term in Eq. 71 reads

$$(\text{FK}) R D^{-1/2} e^{-\tilde{H}t} \tilde{f}_0 = (\text{FK}) R D^{-1/2} e^{-\tilde{H}t} D^{1/2} R^T f_0 \tag{78}$$

$$= \sum_k \frac{(-t)^k}{k!} (\text{FK}) (\Theta M)^k f_0 \tag{79}$$

$$= \sum_k \frac{(-t)^k}{k!} [(\text{FK}) \Theta (\text{FK})^T C_Y^{-1}]^k (\text{FK}) f_0 , \tag{80}$$

and hence the contribution to $\tilde{\epsilon}_t$

$$D_Y^{1/2} R_Y^T R_Y D_Y^{-1/2} e^{-\tilde{H}_\epsilon t} D_Y^{1/2} R_Y^T (\text{FK}) f_0 = e^{-\tilde{H}_\epsilon t} D_Y^{1/2} R_Y^T ((\text{FK}) f_0) . \tag{81}$$

Combining Eqs. 77 and 81 yields

$$e^{-\tilde{H}_\epsilon t} D_Y^{1/2} R_Y^T (y - (\text{FK}) f_0) - \tilde{\epsilon}_\infty = e^{-\tilde{H}_\epsilon t} \tilde{\epsilon}_0 - \tilde{\epsilon}_\infty . \tag{82}$$

And finally, adding Eqs. 69 and 82, we obtain

$$\tilde{\epsilon}_t = e^{-\tilde{H}_\epsilon t} \tilde{\epsilon}_0 , \tag{83}$$

as expected.

2.1 Evolution reloaded

A basis in the N_{grid} -dimensional space of vectors $f_t \in \mathbb{R}^{N_{\text{grid}}}$, $\{z^{(k)}; k = 1, \dots, N_{\text{grid}}\}$, can be constructed from the eigenvectors of Θ ,

$$\Theta z^{(k)} = \lambda^{(k)} z^{(k)} . \tag{84}$$

Since Θ is symmetric and positive semi-definite, we know that $\lambda^{(k)} \geq 0$ and that the eigenvectors are a complete set of orthonormal vectors. It is useful to introduce the following notation:

$$k \in \mathcal{I}_{\Theta^\perp} , \text{ if } \lambda^{(k)} \neq 0 , \quad k \in \mathcal{I}_{\Theta^\parallel} , \text{ if } \lambda^{(k)} = 0 . \tag{85}$$

The dimension of the kernel of Θ is $N_{\Theta\parallel} = |\mathcal{I}_{\Theta\parallel}|$, the dimension of its orthogonal complement is $N_{\Theta\perp} = |\mathcal{I}_{\Theta\perp}|$, and $N_{\Theta\parallel} + N_{\Theta\perp} = N_{\text{grid}}$.

Projecting the flow equation onto the basis vectors $z^{(k)}$ yields

$$\frac{d}{dt}f_{t,k} = 0, \quad \text{for } k \in \mathcal{I}_{\Theta\parallel}, \quad (86)$$

$$\frac{d}{dt}f_{t,k} = -\lambda^{(k)} \left(z^{(k)}, Mf_t \right) + \lambda^{(k)} \left(z^{(k)}, (\text{FK})^T C_Y^{-1} y \right), \quad \text{for } k \in \mathcal{I}_{\Theta\perp}, \quad (87)$$

where $f_{t,k} = (z^{(k)}, f_t)$. Using the completeness of the basis,

$$\left(z^{(k)}, Mf_t \right) = \sum_{k'} \left(z^{(k)}, Mz^{(k')} \right) f_{t,k'} \quad (88)$$

$$= \sum_{k' \in \mathcal{I}_{\Theta\perp}} \left(z^{(k)}, Mz^{(k')} \right) f_{t,k'} + \sum_{k' \in \mathcal{I}_{\Theta\parallel}} \left(z^{(k)}, Mz^{(k')} \right) f_{0,k'}. \quad (89)$$

Finally, we can rewrite Eq. (87) as

$$\frac{d}{dt}f_{t,k} = - \sum_{k' \in \mathcal{I}_{\Theta\perp}} H_{kk'}^{\perp} f_{t,k'} + B_k, \quad (90)$$

where

$$H_{kk'}^{\perp} = \lambda^{(k)} \left(z^{(k)}, Mz^{(k')} \right), \quad (91)$$

is a $N_{\Theta\perp} \times N_{\Theta\perp}$ matrix, and

$$B_k = - \sum_{k' \in \mathcal{I}_{\Theta\parallel}} \lambda^{(k)} \left(z^{(k)}, Mz^{(k')} \right) f_{0,k'} + \lambda^{(k)} \left(z^{(k)}, (\text{FK})^T C_Y^{-1} y \right), \quad (92)$$

is a $N_{\Theta\perp}$ -dimensional vector, with $k \in \mathcal{I}_{\Theta\perp}$.

The other basis, $\{v^{(m)}; m = 1, \dots, N_{\text{grid}}\}$, is made of eigenvectors of M ,

$$Mv^{(m)} = \sigma^{(m)}v^{(m)}. \quad (93)$$

Similarly to the convention introduced above, the eigenvectors in the kernel of M are denoted by indices in $\mathcal{I}_{M\parallel}$, while the eigenvectors in the orthogonal complement are labelled by indices in $\mathcal{I}_{M\perp}$.

2.2 Evolution of f_t revisited.

We start again from the eigenvectors and eigenvalues of Θ

$$\Theta e^{(k)} = \lambda_k e^{(k)}, \quad (94)$$

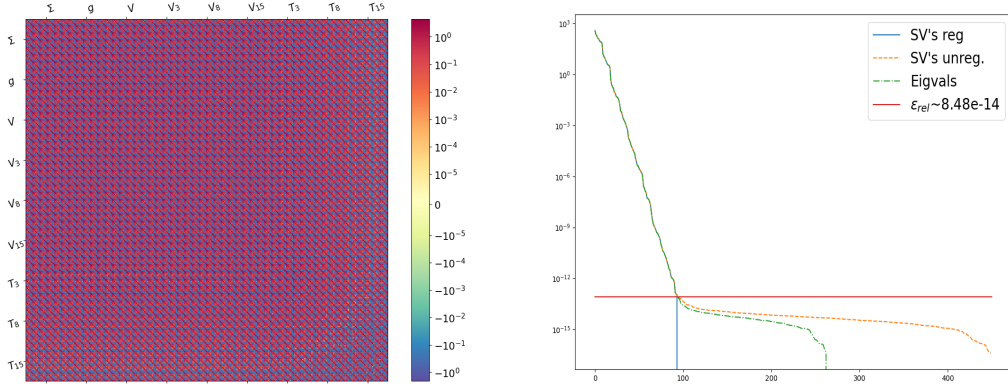


Figure 3: (Left) Representation of the NTK matrix in logarithmic scale. (Right) Decomposition of the NTK matrix in terms of the eigenvalues and of the singular values of the matrix M , and the subsequent regularisation (see main text).

to which we apply the same analysis carried out in Appendix B. In this section, we use a feed-forward neural network with architecture $[1, 28, 20, 9]$. The activation function used in both the two deep layers is `tanh`, while the output function is linear. The weights, rescaled by the square root of the size of the previous layer, are initialised using a normal distribution. The biases are initialised to zero. The input data are the points of the grid in x -space taken from the FK tables.

We start by looking at the NTK matrix, displayed in the left panel of Fig. 3. Upon inspecting the figure, we see that the NTK is uniform across flavours and bins in x . This is indeed expected, as the NTK does not bring any physical information and it does not favour any particular flavour to start with. In other words, the left plot in Fig. 3 shows the amount of information per flavour that we can ever reach given the network architecture and the choice of the grid in x .

The right panel of Fig. 3 shows the decomposition of the NTK matrix in terms of the singular values and in terms of the eigenvalues. The relative machine precision for double precision floating numbers ($\epsilon \sim 2.2 \times 10^{-16}$) is also shown. We see that the singular values and eigenvalues start to diverge when machine precision is reached. This is a symptom of noise in the decomposition, and values below this threshold are set to zero. Above threshold, the singular values and eigenvalues are consistent and they show the usual hierarchy of the NTK extensively discussed in the literature.

The physical information is brought in by the FK tables, which combine with NTK in the flow equation, reported below for convenience

$$\frac{d}{dt} \mathbf{f}_t = \mathbf{b} - \Theta M \mathbf{f}_t. \quad (95)$$

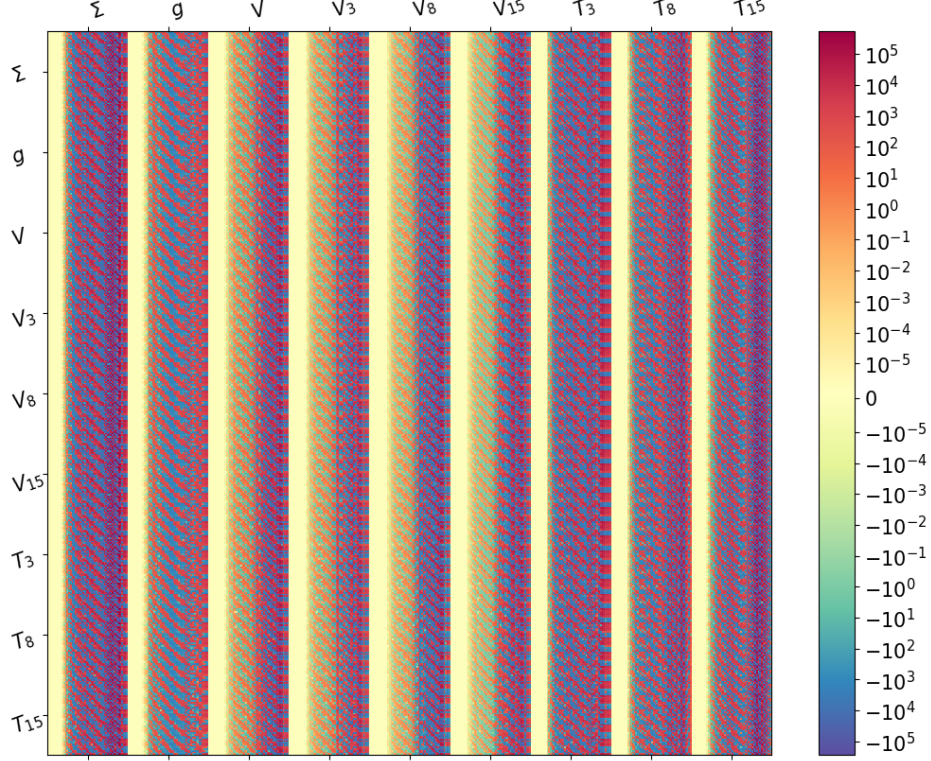


Figure 4:

In a standard ML application of neural networks, the evolution kernel in eq. 95 would be the NTK solely. On the other hand, in the context of PDF fitting, we are interested in solving the inverse problem that arise from the factorization theorem, and that is encoded in the matrix M . This introduces additional complications, as the inverse problem is generally ill-posed and the matrix M is not well-defined (see Appendix B). Moreover, the product taking place in the second term on the rhs of the equation above mixes the NTK and the FK spaces, making it harder to disentangle the information brought in by the NTK and the FK tables separately.

The resulting evolution kernel is shown in Fig. 4. Compared to the left panel of Fig. (3), the kernel preserves the same underlying structure, yet modified by the action of the FK tables. We also see that the overall architecture (network and FK tables) is not sensitive to the low- x region, while the mid- and large- x region is more sensitive.

Unfortunately, the evolution kernel is singular and cannot be inverted as it is, preventing any attempt to integrate the flow equation eq. (95). In order to proceed, we need to decompose the kernel in its various components and discard the zero modes. This is the subject of the next section.

Projection of the flow equation

The evolution kernel cannot be inverted directly, due to the presence of many zero modes. We then resort to working in a subspace of the kernel, where it can be inverted. In doing so, we will also try to disentangle the effects of the NTK and of the FK tables.

We start by projecting the flow equation in the eigenspace of the NTK, so that we have

$$\frac{d}{dt}f'_{k,t} = b'_k - \lambda_k(\mathbf{e}^{(k)}, M\mathbf{f}_t). \quad (96)$$

Here the prime denotes the k -th component of the vector in the eigenspace of the NTK, and $(\mathbf{e}^{(k)}, M\mathbf{f}_t)$ is the ordinary scalar product. Eq. (96) can also be recast into matrix form. Writing $\Theta = U_\Theta \Sigma_\Theta U_\Theta^T$, then we have

$$\frac{d}{dt}\mathbf{f}'_t = \mathbf{b}' - \Sigma_\Theta U_\Theta^T M U_\Theta \mathbf{f}' = \mathbf{b}' - \Sigma_\Theta M' \mathbf{f}'. \quad (97)$$

The equation above requires two observations. The first one is that some eigenvalues of the NTK are actually zero, so that the second term on the rhs of eq. (96) is zero. Furthermore, for the first term on the rhs we have

$$b'_k = (\mathbf{e}^{(k)}, \mathbf{b}) = (\mathbf{e}^{(k)}, \Theta(FK)^T C_Y^{-1} \mathbf{Y}) = 0 \quad \forall \mathbf{e}^{(k)} \in \text{Ker}(\Theta). \quad (98)$$

Thus, for zero eigenvalues of the NTK, we have

$$\frac{d}{dt}f'_{k,t} = 0 \quad \Rightarrow \quad f'_{k,t} = f'_{k,0}, \quad (99)$$

and this result holds regardless of the FK tables.

The second observation involves the matrix M . In particular, we can decompose M as follows

$$M = LL^T = (RD^{1/2})(D^{1/2}R^T), \quad (100)$$

where D is the diagonal matrix of the eigenvalues (or singular values) of M and R is the matrix of the eigenvectors. Using the completeness relation of the vectors $\mathbf{e}^{(k)}$, we can write

$$M'_{k_1 k_2} = (e^{(k_1)}, L e^{(k)}) (e^{(k)}, L^T e^{(k_2)}) = L'_{k_1 k} (L^T)'_{k k_2} = L'_{k_1 k} L'_{k_2 k}, \quad (101)$$

or, in matrix form,

$$M' = U_\Theta^T M U_\Theta = U_\Theta^T L U_\Theta U_\Theta^T L^T U_\Theta = L'(L')^T, \quad (102)$$

where in the second equality we have used the completeness relation of the eigenvectors.

A close-up inspection of L' reveals that

$$\begin{aligned}
L'_{k_1 k_2} &\equiv \left(\mathbf{e}^{(k)}, RD^{1/2} \mathbf{e}^{(k)} \right) \\
&= e_{i_1}^{(k_1)} R_{i_1 i_2} \mu_{i_2} e_{i_2}^{(k_2)} \\
&= e_{i_1}^{(k_1)} v_{i_1}^{(i_2)} \mu_{i_2} e_{i_2}^{(k_2)} \\
&= \mu_{i_2} e_{i_2}^{(k_2)} (\mathbf{e}^{(k_1)}, \mathbf{v}^{(i_2)}),
\end{aligned} \tag{103}$$

which means that the intersection between the NTK and the FK tables takes place in the scalar product of the last equality above.

By means of eq. (101), we can rewrite the flow equation as follows

$$\frac{d}{dt} f'_{k_1, t} = b'_{k_1} - \lambda_{k_1} L'_{k_1 k_2} L'_{k_3 k_2} f'_{k_3, t}. \tag{104}$$

Finally, in order to obtain a symmetric kernel, we can multiply both sides by $(L')^T$ so that we have

$$\frac{d}{dt} \tilde{f}'_{k_1, t} = \tilde{b}'_{k_1} - \tilde{K}_{k_1 k_2} \tilde{f}'_{k_2, t} \tag{105}$$

where $\tilde{f}'_{k_1, t}$ and \tilde{b}'_{k_1} are the components rotated by the action of L' , and $\tilde{K}_{k_1 k_2}$ is the symmetric kernel defined as

$$\tilde{K} = (L')^T \Sigma_{\Theta} L'. \tag{106}$$

Removing zero modes of the NTK

In obtaining Eq. (105), we have kept all the zero modes of the NTK. However, we can filter them out and write an expression that contains zero modes of the FK tables solely. To do so, we restart from eq. (97), and we recall that the rotations applied to the vectors are

$$\mathbf{b}' = U_{\Theta}^T \mathbf{b} \quad \text{and} \quad \mathbf{f}'_t = U_{\Theta}^T \mathbf{f}_t. \tag{107}$$

In order to separate that zero modes, we can split the vector \mathbf{f}'_t into a part that is orthogonal to the kernel of the NTK, and another which lives in the kernel and hence is constant

$$\frac{d}{dt} \begin{pmatrix} \mathbf{f}'_{\perp, t} \\ \mathbf{f}'_{\parallel, t} \end{pmatrix} = \begin{pmatrix} \mathbf{b}'_{\perp} \\ \mathbf{b}'_{\parallel} \end{pmatrix} - \begin{pmatrix} \Sigma_{\Theta} & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} M'_{\perp\perp} & M'_{\perp\parallel} \\ M'_{\parallel\perp} & M'_{\parallel\parallel} \end{pmatrix} \begin{pmatrix} \mathbf{f}'_{\perp, t} \\ \mathbf{f}'_{\parallel, t} \end{pmatrix} \tag{108}$$

where the subscripts \perp and \parallel denote the components orthogonal and parallel to the kernel of the NTK, respectively. The equation above unravels as follows

$$\frac{d}{dt} \mathbf{f}'_{\perp, t} = \mathbf{b}'_{\perp} - \Sigma_{\Theta} \left(M'_{\perp\perp} \mathbf{f}'_{\perp, t} + M'_{\perp\parallel} \mathbf{f}'_{\parallel, t} \right), \tag{109}$$

$$\frac{d}{dt} \mathbf{f}'_{\parallel, t} = \mathbf{0}, \tag{110}$$

where the second equation is a consequence of the fact that the kernel of the NTK is a zero mode for the learning flow. Note that the parallel component $\mathbf{f}'_{\parallel,t}$ appears in the flow equation of $\mathbf{f}'_{\perp,t}$, but it does not evolve in time. Hence, it can be reabsorbed in the constant factor \mathbf{b}'_{\perp} , giving

$$\frac{d}{dt}\mathbf{f}'_{\perp,t} = \mathbf{B}'_{\perp} - \Sigma_{\Theta}M'_{\perp\perp}\mathbf{f}'_{\perp,t}, \quad (111)$$

where

$$\mathbf{B}'_{\perp} = \mathbf{b}'_{\perp} - \Sigma_{\Theta}M'_{\perp\parallel}\mathbf{f}'_{\parallel,t}. \quad (112)$$

Eq. (111) restricts the evolution of the vector \mathbf{f}'_t to the components orthogonal to the kernel of the NTK. The only flat directions left are those encoded in the FK tables.

Flat directions of the FK tables

Following the same reasoning that led to eq. (105), we can rewrite eq. (111) as follows

$$\frac{d}{dt}\tilde{\mathbf{f}}'_{\perp,t} = \tilde{\mathbf{B}}'_{\perp} - \tilde{K}_{\perp}\tilde{\mathbf{f}}'_{\perp,t}, \quad (113)$$

where

$$M'_{\perp\perp} = L'_{\perp}(L'_{\perp})^T, \quad \tilde{\mathbf{f}}'_{\perp,t} = (L'_{\perp})^T\mathbf{f}'_{\perp,t}, \quad \text{and} \quad \tilde{K}_{\perp} = (L'_{\perp})^T\Sigma_{\Theta}L'_{\perp}. \quad (114)$$

The entries of the symmetric kernel \tilde{K}_{\perp} and its decomposition are shown in the left and right panel of Fig. ... As previously done for then NTK, we can rewrite this kernel as follows

$$\tilde{K}_{\perp} = U_K\Sigma_KU_K^T, \quad (115)$$

and rotate equation eq. (113) to the eigenbasis of the kernel

$$\frac{d}{dt}\hat{\mathbf{f}}_{\perp} = \hat{\mathbf{B}}_{\perp} - \Sigma_K\hat{\mathbf{f}}_{\perp,t}, \quad (116)$$

where

$$\hat{\mathbf{f}}_{\perp} = U_K^T\tilde{\mathbf{f}}'_{\perp}, \quad \text{and} \quad \hat{\mathbf{B}}_{\perp} = U_K^T\tilde{\mathbf{B}}'_{\perp}. \quad (117)$$

As displayed in Fig. ..., the kernel contains many flat directions due to the FK tables. We can then decompose the vectors as follows

$$\frac{d}{dt}\begin{pmatrix} \hat{\mathbf{f}}_{\perp\perp,t} \\ \hat{\mathbf{f}}_{\perp\parallel,t} \end{pmatrix} = \begin{pmatrix} \hat{\mathbf{B}}_{\perp\perp} \\ \hat{\mathbf{B}}_{\perp\parallel} \end{pmatrix} - \begin{pmatrix} \Sigma_K & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \hat{\mathbf{f}}_{\perp\perp,t} \\ \hat{\mathbf{f}}_{\perp\parallel,t} \end{pmatrix} \quad (118)$$

yielding the two independent equations

$$\frac{d}{dt}\hat{\mathbf{f}}_{\perp\perp,t} = \hat{\mathbf{B}}_{\perp\perp} - \Sigma_K\hat{\mathbf{f}}_{\perp\perp,t} \quad (119)$$

$$\frac{d}{dt}\hat{\mathbf{f}}_{\perp\parallel,t} = 0. \quad (120)$$

Once inverted, the first equation provides the evolution of the components of the vector orthogonal to the null space of both NTK and FK tables. To retrieve the original vector in flavour basis, we need to apply to the solution for $\hat{\mathbf{f}}_{\perp\perp,t}$ the inverse rotations applied so far. This gives

$$\mathbf{f}_t = [U_\Theta]_\perp L'_\perp U_K \hat{\mathbf{f}}_{\perp\perp,t}, \quad (121)$$

where $[U_\Theta]_\perp$ is the subspace of the eigenvectors of the NTK that is orthogonal to its kernel.

Properties of the symmetric kernel

The kernel \tilde{K} ...

Obsolete Material – to be deleted, later

We then project the flow equation for f_t , Eq. 50, in this basis

$$\frac{d}{dt}f_{i,t} = b_i - \lambda_i M_{ij} f_{j,t}, \quad (122)$$

where

$$f_{i,t} = (v_i, f_t), \quad (123)$$

$$b_i = (v_i, \Theta(\text{FK})^T C_Y^{-1} y) \quad (124)$$

$$M_{ij} = (v_i, (\text{FK})^T C_Y^{-1} (\text{FK}) v_j). \quad (125)$$

Note that M is symmetric and positive definite. We can write $M = R D^{1/2} D^{1/2} R^T$ and repeat the derivation we did for the ϵ .

To be continued...

At order $O(1)$ in an expansion in $1/n$, where n is the width of the NN, Θ is constant during training, and the flow equation can be integrated analytically,

$$f_t = M^{-1} (\text{FK})^T \left[\left(1 - e^{-\hat{\Theta}_Y t}\right) y + e^{-\hat{\Theta}_Y t} (\text{FK}) f_0 \right]. \quad (126)$$

where

$$\Theta_Y = (\text{FK})^T \Theta (\text{FK}), \quad (127)$$

$$\hat{\Theta}_Y = \Theta_Y C_Y^{-1}, \quad (128)$$

$$M = (\text{FK})^T (\text{FK}). \quad (129)$$

Note that, if (FK) were invertible, then for $t \rightarrow \infty$

$$(\text{FK}) f_\infty = (\text{FK}) M^{-1} (\text{FK})^T y = y, \quad (130)$$

and the theoretical predictions go exactly through the points. However this is not true in the real life scenario where (FK) is *not* invertible. It is therefore interesting to compute the matrix (in data space)

$$(\text{FK}) M^{-1} (\text{FK})^T, \quad (131)$$

in order to understand how close the training can get to the training points.

For points f^* that do not enter in the computation of the theory prediction, the flow equation is

$$\frac{d}{dt} f_t^* = \Theta_t^* (\text{FK})^T C_Y^{-1} \epsilon_t, \quad (132)$$

where

$$\Theta_t^* = \sum_{\mu, \nu} \lambda_{\mu\nu} (\nabla_\mu f_t^*) (\nabla_\nu f_t)^T. \quad (133)$$

The solution to this flow equation is

$$f_t^* = \Theta^* \Theta^{-1} f_t. \quad (134)$$

Need to check the constant term so that the boundary condition at $t = 0$ is correct.

References

- [1] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and generalization in neural networks. *CoRR*, abs/1806.07572, 2018.

A Null space of FK

The FK tables can be regarded as a linear map from the space of PDFs to the space of data:

$$(\text{FK}) : \mathbb{R}^{\text{PDF}} \rightarrow \mathbb{R}^{N_{\text{dat}}} . \quad (135)$$

Note that the matrix corresponding to this linear map is not square, and hence $(\text{FK}) \neq (\text{FK})^T$. We can define the null space of the FK tables

$$\ker(\text{FK}) \equiv K_{(\text{FK})} = \{f \in \mathbb{R}^{\text{PDF}} : (\text{FK})f = 0\} , \quad (136)$$

together with its orthogonal space

$$R_{(\text{FK})} \equiv K_{(\text{FK})}^\perp = \{f \in \mathbb{R}^{\text{PDF}} : f \cdot f_K = 0, \forall f_K \in K_{(\text{FK})}\} . \quad (137)$$

Note that

$$(\text{FK})f = 0 \implies Mf = 0 . \quad (138)$$

The converse is also true,

$$Mf = 0 \implies (f, Mf) = 0 \quad (139)$$

$$\implies ((\text{FK})f, C_Y^{-1}(\text{FK})f) = 0 \quad (140)$$

$$\implies (\text{FK})f = 0 , \quad (141)$$

where the last inequality follows from the fact that $C_Y > 0$.

For each of these two subspaces we can define a basis

$$\mathcal{B}_K = \left\{ \mathbf{u}_K^{(i)} , \ i = 1, \dots, \dim K_{(\text{FK})} \right\} , \quad (142)$$

$$\mathcal{B}_\perp = \left\{ \mathbf{u}_\perp^{(i)} , \ i = 1, \dots, \dim R_{(\text{FK})} \right\} . \quad (143)$$

Remember that $K_{(\text{FK})} \oplus R_{(\text{FK})} = \mathbb{R}^{\text{PDF}}$ and thus the basis $\mathcal{B} = \mathcal{B}_K \oplus \mathcal{B}_\perp$ is a basis for \mathbb{R}^{PDF} . Henceforth, when decomposing a vector in \mathbb{R}^{PDF} , I will use the ordering $\{\mathcal{B}_K, \mathcal{B}_\perp\}$. Hence, given $\mathbf{f} \in \mathbb{R}^{\text{PDF}}$, we can write

$$\mathbf{f} = \mathbf{f}_K + \mathbf{f}_\perp = \begin{pmatrix} f_K \\ f_\perp \end{pmatrix} . \quad (144)$$

Finally, note that (FK) is not symmetric (not even square). Thus, the right null space is not the same as the left null space, in particular

$$(\text{FK})\mathbf{u}_K = 0 \not\Rightarrow \mathbf{u}_K^T(\text{FK}) = 0 . \quad (145)$$

With the two bases in eqs. (143), we can decompose the (FK) as follows

$$(\text{FK}) = \begin{pmatrix} (\text{FK})_{KK} & (\text{FK})_{K\perp} \\ (\text{FK})_{\perp K} & (\text{FK})_{\perp\perp} \end{pmatrix}, \quad (146)$$

where

$$(\text{FK})_{B_1 B_2} = \sum_{i=1}^{\dim \mathcal{B}_1} \sum_{j=1}^{\dim \mathcal{B}_2} \langle u_{\mathcal{B}_1}^{(i)} | (\text{FK}) | u_{\mathcal{B}_2}^{(j)} \rangle \quad \mathcal{B}_1, \mathcal{B}_2 = K_{(\text{FK})}, R_{(\text{FK})}. \quad (147)$$

By definition of the kernel, we also have

$$(\text{FK})_{KK} = (\text{FK})_{\perp K} = 0, \quad (148)$$

while $(\text{FK})_{K\perp}$ would be zero only if (FK) was diagonal. Thus, in the basis \mathcal{B} , the FK tables can be expressed as

$$(\text{FK}) = \begin{pmatrix} 0 & (\text{FK})_{K\perp} \\ 0 & (\text{FK})_{\perp\perp} \end{pmatrix}. \quad (149)$$

The product with a vector $\mathbf{f} \in \mathbb{R}^{\text{PDF}}$ becomes

$$(\text{FK})\mathbf{f} = \begin{pmatrix} 0 & (\text{FK})_{K\perp} \\ 0 & (\text{FK})_{\perp\perp} \end{pmatrix} \begin{pmatrix} f_K \\ f_\perp \end{pmatrix} = \begin{pmatrix} (\text{FK})_{K\perp} \\ (\text{FK})_{\perp\perp} \end{pmatrix} f_\perp. \quad (150)$$

The matrix M , eq. (46), can also be decomposed as in Eq. (147). However, now $M = M^T$ and thus only one component is non-zero:

$$M = \begin{pmatrix} 0 & 0 \\ 0 & M_{\perp\perp} \end{pmatrix}. \quad (151)$$

We can decompose the other elements in eq. (52) in a similar vein:

$$b = \begin{pmatrix} b_K \\ b_\perp \end{pmatrix}, \quad \Theta = \begin{pmatrix} \Theta_{KK} & \Theta_{K\perp} \\ \Theta_{\perp K} & \Theta_{\perp\perp} \end{pmatrix}, \quad (152)$$

where $\Theta \in \mathbb{R}^{\text{PDF} \times \text{PDF}}$ and can thus be projected in \mathcal{B} . Owing to this decomposition, we can now split eq. (52) into the two components of f_t , so that we are left with two coupled linear differential equations:

$$\begin{cases} \frac{d}{dt} f_\perp = b_\perp - \Theta_{\perp\perp} M_{\perp\perp} f_\perp, \\ \frac{d}{dt} f_K = b_K - \Theta_{K\perp} M_{\perp\perp} f_\perp, \end{cases} \quad (153)$$

where the dependence on the training time t is left implicit. From eq. (153) we see that f_K does *not* evolve linearly, but depends on the evolution of f_\perp . Moreover, the differential equation for f_\perp is diagonal, and can be solved as shown in the main text, namely diagonalising $M_{\perp\perp}$. Note that, in the orthogonal subspace, $M_{\perp\perp}$ is positive-definite as we have projected out the null eigenvalues.

B Numerical validation of eq. (67)

In eq. (67) we showed that the solution for ε_∞ in the limit of infinite training length corresponds to the zero mode of the evolution matrix in data space \tilde{H}_ε . Here we provide the numerical argument that supports this result.

We start from the matrix M which, roughly speaking, is the contraction of two FK tables. Since FK tables contain many zeros, the matrix M will be sparse. The matrix M is displayed element-wise in Fig. 5. Each block identifies a flavour, and the length of the side of the block is equal to the size of the grid. The matrix M is decomposed using the singular

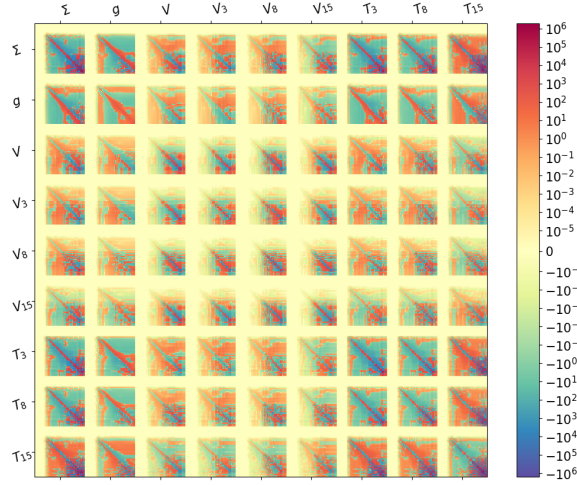


Figure 5: Matrix M in logarithmic scale.

value decomposition (SVD) method

$$M = U\Sigma V^T, \quad (154)$$

where U and V are orthogonal matrices, and Σ is a diagonal matrix containing the singular values sorted in decreasing order. Recall that the matrix M is also symmetric by construction, and this guarantees the existence of the solution of the associated eigensystem. Singular values and the absolute eigenvalues are displayed in Fig. 6. Here we also show the relative machine precision error, defined as the product of the largest singular value with the machine precision for double precision floating numbers ($\epsilon \sim 2.2 \times 10^{-16}$). Since the largest singular value is of order 10^6 , the relative machine precision error is $\epsilon_{\text{rel}} \sim 1.2 \times 10^{-9}$. The relevance of this number will be clear in the following observation.

As it can be seen from the figure, the singular values and the eigenvalues overlap as long as the relative numerical precision is not reached. Once this numerical threshold is exceeded, the two sets of values start to diverge. This is a consequence of the fact that ϵ_{rel}

defines the smallest double floating number that machine precision allows us to distinguish. Everything below this value is noise and should be discarded. This prompts us to apply an *ad hoc* regularisation to Σ , whereby the values below ϵ_{rel} are set to zero. Note that we have

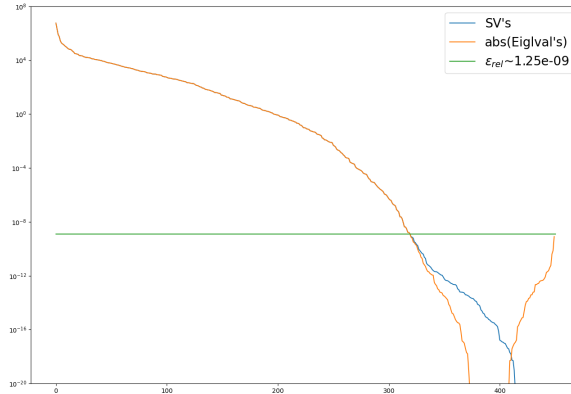


Figure 6: Singular values and (absolute) eigenvalues of the matrix M plotted in logarithmic scale. The green horizontal line marks the relative machine precision error, as explained in the main text.

also checked that the singular vectors of U and V are the same if they are associated to the singular values above the numerical threshold. In other words, in the space orthogonal to $\text{Ker}(M)$, the SVD matches the eigendecomposition.

In principle, the singular vectors of the non-zero values span the space orthogonal to $\text{Ker}(M)$. In this subspace, the matrix M is diagonal and the entries are exactly the non-zero singular values. However, if we compute the condition number of M in this subspace, we find that it is remarkably large, $\kappa(M) \sim 10^{15}$. In other words, even in this subspace the matrix M is ill-defined and cannot be inverted. The solution to this problem is to further reduce the dimensionality of the orthogonal space, at the cost of introducing some arbitrariness into the problem. Specifically, the final solution will depend on the number of small singular values that we discard. For example, the solution at infinity, eq. (58), will be influenced by this cut-off, as shown in Fig. 7. In this figure, eq. (58) is first computed in the truncated orthogonal space; then the result is transformed back to the original flavour basis using the matrix U restricted to the vectors that span the (reduced) orthogonal space.

As a final check, we can verify that eq. (67) is numerically satisfied. As we have just mentioned, even ϵ_∞ will depend on the severity of the cut-off. We then expect different results as we modify the number of singular values that we discard. For instance, using the orthogonal subspace without any cut, eq. (67) evaluates to

$$\left| \tilde{H}_\epsilon \epsilon_\infty \right| \approx 0.28. \quad (155)$$

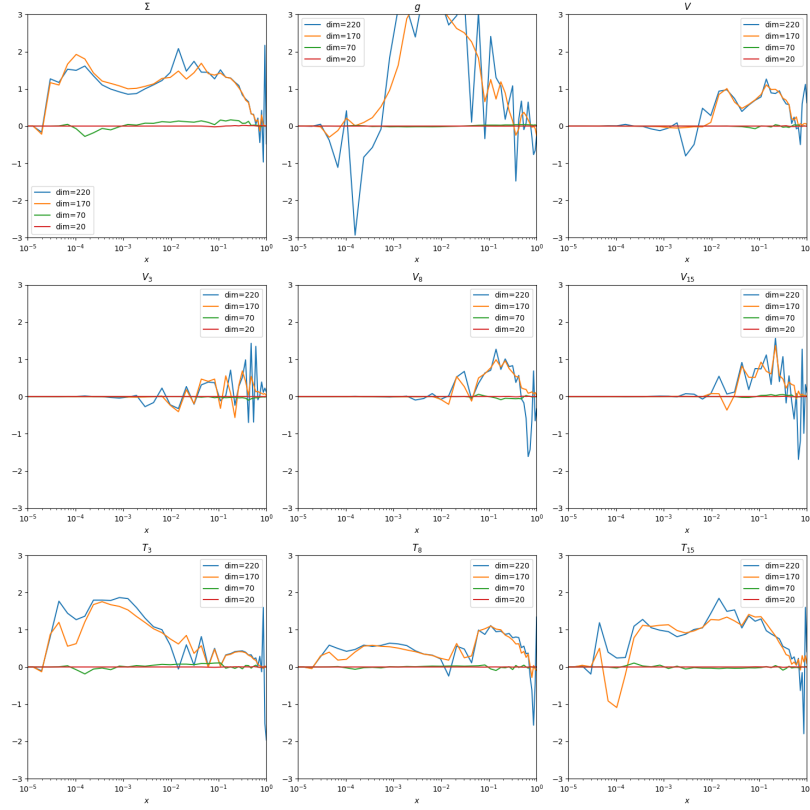


Figure 7:

On the other hand, if we reduce the size of the subspace to 170 (from 320), we get

$$\left| \tilde{H}_\varepsilon \varepsilon_\infty \right| \approx 1.33 \times 10^{-5}. \quad (156)$$