

---

# Self-Distillation Amplifies Regularization in Hilbert Space Supplementary Appendix

---

Hossein Mobahi<sup>♣</sup> Mehrdad Farajtabar<sup>§</sup> Peter L. Bartlett<sup>♣‡</sup>

[hmobahi@google.com](mailto:hmobahi@google.com) [farajtabar@google.com](mailto:farajtabar@google.com) [bartlett@eecs.berkeley.edu](mailto:bartlett@eecs.berkeley.edu)

♣ Google Research, Mountain View, CA, USA

§ DeepMind, Mountain View, CA, USA

‡ EECS Dept., University of California at Berkeley, Berkeley, CA, USA

## A Solving the Variational Problem

In this section we derive the solution to the following variational problem,

$$f^* \triangleq \arg \min_{f \in \mathcal{F}} \frac{1}{K} \sum_k \left( f(\mathbf{x}_k) - y_k \right)^2 + c \int_{\mathcal{X}} \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f(\mathbf{x}) f(\mathbf{x}^\dagger) d\mathbf{x} d\mathbf{x}^\dagger. \quad (26)$$

Using Dirac delta function, we can rewrite the objective function as,

$$f^* = \arg \min_{f \in \mathcal{F}} \frac{1}{K} \sum_k \left( \int_{\mathcal{X}} f(\mathbf{x}) \delta(\mathbf{x} - \mathbf{x}_k) d\mathbf{x} - y_k \right)^2 + c \int_{\mathcal{X}} \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f(\mathbf{x}) f(\mathbf{x}^\dagger) d\mathbf{x} d\mathbf{x}^\dagger. \quad (27)$$

For brevity, name the objective functional  $J$ ,

$$J(f) \triangleq \frac{1}{K} \sum_k \left( \int_{\mathcal{X}} f(\mathbf{x}) \delta(\mathbf{x} - \mathbf{x}_k) d\mathbf{x} - y_k \right)^2 + c \int_{\mathcal{X}} \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f(\mathbf{x}) f(\mathbf{x}^\dagger) d\mathbf{x} d\mathbf{x}^\dagger. \quad (28)$$

If  $f^*$  minimizes the  $J(f)$ , it must be a stationary point of  $J$ . That is,  $J(f + \epsilon\phi) = J(f)$ , for any  $\phi \in \mathcal{F}$  as  $\epsilon \rightarrow 0$ . More precisely, it is necessary for  $f^*$  to satisfy,

$$\forall \phi \in \mathcal{F}; \left( \frac{d}{d\epsilon} J(f^* + \epsilon\phi) \right)_{\epsilon=0} = 0. \quad (29)$$

We first construct  $J(f^* + \epsilon\phi)$ ,

$$J(f^* + \epsilon\phi) = \frac{1}{K} \sum_k \left( \int_{\mathcal{X}} [f^* + \epsilon\phi](\mathbf{x}) \delta(\mathbf{x} - \mathbf{x}_k) d\mathbf{x} - y_k \right)^2 \quad (30)$$

$$+ c \int_{\mathcal{X}} \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) [f^* + \epsilon\phi](\mathbf{x}) [f^* + \epsilon\phi](\mathbf{x}^\dagger) d\mathbf{x} d\mathbf{x}^\dagger, \quad (31)$$

or equivalently,

$$J(f^* + \epsilon\phi) = \frac{1}{K} \sum_k \left( \int_{\mathcal{X}} (f^*(\mathbf{x}) + \epsilon\phi(\mathbf{x})) \delta(\mathbf{x} - \mathbf{x}_k) d\mathbf{x} - y_k \right)^2 \quad (32)$$

$$+ c \int_{\mathcal{X}} \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) (f^*(\mathbf{x}) + \epsilon\phi(\mathbf{x})) (f^*(\mathbf{x}^\dagger) + \epsilon\phi(\mathbf{x}^\dagger)) d\mathbf{x} d\mathbf{x}^\dagger. \quad (33)$$

Thus,

$$\frac{d}{d\epsilon} J(f^* + \epsilon\phi) \quad (34)$$

$$= \frac{1}{K} \sum_k 2 \left( \int_{\mathcal{X}} (f^*(\mathbf{x}^\diamond) + \epsilon\phi(\mathbf{x}^\diamond)) \delta(\mathbf{x}^\diamond - \mathbf{x}_k) d\mathbf{x}^\diamond - y_k \right) \left( \int_{\mathcal{X}} \phi(\mathbf{x}) \delta(\mathbf{x} - \mathbf{x}_k) d\mathbf{x} \right) \quad (35)$$

$$+ c \int_{\mathcal{X}} \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) \left( \phi(\mathbf{x}) (f^*(\mathbf{x}^\dagger) + \epsilon\phi(\mathbf{x}^\dagger)) + \phi(\mathbf{x}^\dagger) (f^*(\mathbf{x}) + \epsilon\phi(\mathbf{x})) \right) d\mathbf{x} d\mathbf{x}^\dagger. \quad (36)$$

Setting  $\epsilon = 0$ ,

$$\left( \frac{d}{d\epsilon} J(f^* + \epsilon\phi) \right)_{\epsilon=0} = \frac{1}{K} \sum_k 2 \left( \int_{\mathcal{X}} f^*(\mathbf{x}^\diamond) \delta(\mathbf{x}^\diamond - \mathbf{x}_k) d\mathbf{x}^\diamond - y_k \right) \left( \int_{\mathcal{X}} \phi(\mathbf{x}) \delta(\mathbf{x} - \mathbf{x}_k) d\mathbf{x} \right) \quad (37)$$

$$+ c \int_{\mathcal{X}} \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) \left( \phi(\mathbf{x}) f^*(\mathbf{x}^\dagger) + \phi(\mathbf{x}^\dagger) f^*(\mathbf{x}) \right) d\mathbf{x} d\mathbf{x}^\dagger. \quad (38)$$

By the symmetry of  $u$ ,

$$\left( \frac{d}{d\epsilon} J(f^* + \epsilon\phi) \right)_{\epsilon=0} = \frac{1}{K} \sum_k 2 \left( \int_{\mathcal{X}} f^*(\mathbf{x}^\diamond) \delta(\mathbf{x}^\diamond - \mathbf{x}_k) d\mathbf{x}^\diamond - y_k \right) \left( \int_{\mathcal{X}} \phi(\mathbf{x}) \delta(\mathbf{x} - \mathbf{x}_k) d\mathbf{x} \right) \quad (39)$$

$$+ 2c \int_{\mathcal{X}} \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) \phi(\mathbf{x}) f^*(\mathbf{x}^\dagger) d\mathbf{x} d\mathbf{x}^\dagger. \quad (40)$$

Factoring out  $\phi$ ,

$$\left( \frac{d}{d\epsilon} J(f^* + \epsilon\phi) \right)_{\epsilon=0} = \int_{\mathcal{X}} 2\phi(\mathbf{x}) \left( \frac{1}{K} \sum_k \delta(\mathbf{x} - \mathbf{x}_k) \left( \int_{\mathcal{X}} f^*(\mathbf{x}^\diamond) \delta(\mathbf{x}^\diamond - \mathbf{x}_k) d\mathbf{x}^\diamond - y_k \right) \right) \quad (41)$$

$$+ c \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f^*(\mathbf{x}^\dagger) d\mathbf{x}^\dagger d\mathbf{x}. \quad (42)$$

In order for the above to be zero for  $\forall \phi \in \mathcal{F}$ , it is necessary that,

$$\frac{1}{K} \sum_k \delta(\mathbf{x} - \mathbf{x}_k) \left( \int_{\mathcal{X}} f^*(\mathbf{x}^\diamond) \delta(\mathbf{x}^\diamond - \mathbf{x}_k) d\mathbf{x}^\diamond - y_k \right) + c \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f^*(\mathbf{x}^\dagger) d\mathbf{x}^\dagger = 0, \quad (43)$$

which further simplifies to,

$$\frac{1}{K} \sum_k \delta(\mathbf{x} - \mathbf{x}_k) (f^*(\mathbf{x}_k) - y_k) + c \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f^*(\mathbf{x}^\dagger) d\mathbf{x}^\dagger = 0. \quad (44)$$

We can equivalently express (44) by the following system of equations,

$$\begin{cases} \frac{1}{K} \sum_k \delta(\mathbf{x} - \mathbf{x}_k) r_k + c \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f^*(\mathbf{x}^\dagger) d\mathbf{x}^\dagger = 0 \\ r_1 = f^*(\mathbf{x}_1) - y_1 \\ \vdots \\ r_K = f^*(\mathbf{x}_K) - y_K \end{cases}. \quad (45)$$

We first focus on solving the first equation in  $f^*$ ,

$$\frac{1}{K} \sum_k \delta(\mathbf{x} - \mathbf{x}_k) r_k + c \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f^*(\mathbf{x}^\dagger) d\mathbf{x}^\dagger = 0; \quad (46)$$

later we can replace the resulted  $f^*$  in other equations to obtain  $r_k$ 's. Let  $g(\mathbf{x}, \mathbf{t})$  be a function such that,

$$\int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) g(\mathbf{x}^\dagger, \mathbf{t}) d\mathbf{x}^\dagger = \delta(\mathbf{x} - \mathbf{t}). \quad (47)$$

Such  $g$  is called the **Green's function** of the linear operator  $L$  satisfying  $[Lf](\mathbf{x}) = \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f(\mathbf{x}^\dagger) d\mathbf{x}^\dagger$ . If we multiply both sides of (47) by  $\frac{1}{K} \sum_k \delta(\mathbf{t} - \mathbf{x}_k) r_k$  and then integrate w.r.t.  $\mathbf{t}$ , we obtain,

$$\int_{\mathcal{X}} \left( \frac{1}{K} \sum_k r_k \delta(\mathbf{t} - \mathbf{x}_k) \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) g(\mathbf{x}^\dagger, \mathbf{t}) d\mathbf{x}^\dagger \right) d\mathbf{t} \quad (48)$$

$$= \int_{\mathcal{X}} \left( \frac{1}{K} \sum_k r_k \delta(\mathbf{t} - \mathbf{x}_k) \delta(\mathbf{x} - \mathbf{t}) \right) d\mathbf{t}. \quad (49)$$

Rearranging the left hand side leads to,

$$\int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) \left( \frac{1}{K} \sum_k \int_{\mathcal{X}} r_k \delta(\mathbf{t} - \mathbf{x}_k) g(\mathbf{x}^\dagger, \mathbf{t}) d\mathbf{t} \right) d\mathbf{x}^\dagger \quad (50)$$

$$= \int_{\mathcal{X}} \left( \frac{1}{K} \sum_k r_k \delta(\mathbf{t} - \mathbf{x}_k) \delta(\mathbf{x} - \mathbf{t}) \right) d\mathbf{t}. \quad (51)$$

Using the sifting property of the delta function this simplifies to,

$$\int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) \left( \frac{1}{K} \sum_k r_k g(\mathbf{x}^\dagger, \mathbf{x}_k) \right) d\mathbf{x}^\dagger = \frac{1}{K} \sum_k r_k \delta(\mathbf{x} - \mathbf{x}_k). \quad (52)$$

We can now use the above identity to eliminate  $\frac{1}{K} \sum_k r_k \delta(\mathbf{x} - \mathbf{x}_k)$  in (46) and thus obtain,

$$\int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) \left( \frac{1}{K} \sum_k r_k g(\mathbf{x}^\dagger, \mathbf{x}_k) \right) d\mathbf{x}^\dagger + c \int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) f^*(\mathbf{x}^\dagger) d\mathbf{x}^\dagger = 0, \quad (53)$$

or equivalently

$$\int_{\mathcal{X}} u(\mathbf{x}, \mathbf{x}^\dagger) \left( \frac{1}{K} \sum_k r_k g(\mathbf{x}^\dagger, \mathbf{x}_k) + c f^*(\mathbf{x}^\dagger) \right) d\mathbf{x}^\dagger = 0. \quad (54)$$

A sufficient (and also necessary, as  $u$  is assumed to have empty null space) for the above to hold is that,

$$f^*(\mathbf{x}) = -\frac{1}{cK} \sum_k r_k g(\mathbf{x}, \mathbf{x}_k). \quad (55)$$

We can now eliminate  $f^*$  in the system of equations (45) and obtain a system that only depends on  $r_k$ 's,

$$\begin{cases} r_1 = -\frac{1}{cK} \sum_k r_k g(\mathbf{x}_1, \mathbf{x}_k) - y_1 \\ \vdots \\ r_K = -\frac{1}{cK} \sum_k r_k g(\mathbf{x}_K, \mathbf{x}_k) - y_K \end{cases}. \quad (56)$$

This is a linear system in  $r_k$  and can be expressed in vector/matrix form,

$$(c\mathbf{I} + \mathbf{G})\mathbf{r} = -c\mathbf{y}. \quad (57)$$

Thus,

$$\mathbf{r} = -c(c\mathbf{I} + \mathbf{G})^{-1}\mathbf{y}, \quad (58)$$

and finally using the definition of  $f^*$  in (55) we obtain,

$$f^*(\mathbf{x}) = -\frac{1}{c}\mathbf{g}_x^T\mathbf{r} = \mathbf{g}_x^T(c\mathbf{I} + \mathbf{G})^{-1}\mathbf{y}. \quad (59)$$

## B Equivalent Kernel Regression Problem

Given a positive definite kernel function  $g(\cdot, \cdot)$ . Recall that the solution of regularized kernel regression after  $t$  rounds of self-distillation has the form,

$$f_t^*(\mathbf{x}) = \mathbf{g}_x^T \mathbf{G}^t \Pi_{i=0}^t (\mathbf{G} + c_i \mathbf{I})^{-1} \mathbf{y}_0. \quad (60)$$

On the other hand, the solution to a standard kernel ridge regression on the same training data with a positive definite kernel  $g^\dagger$  has the form,

$$f^\dagger(\mathbf{x}) = \mathbf{g}_x^T (\mathbf{G}^\dagger + c_0 \mathbf{I})^{-1} \mathbf{y}_0, \quad (61)$$

for which there are standard generalization bounds. We claim  $f_t^*$  can be equivalently written in this standard form by a proper choice of  $g^\dagger$  (as a function of  $g$ ). As a result of that, we show the spectrum of the Gram matrix  $\mathbf{G}^\dagger$  relates to that of  $\mathbf{G}$  via,

$$\lambda_k^\dagger = c_0 \frac{1}{\frac{\Pi_{i=0}^t (\lambda_k + c_i)}{\lambda_k^{t+1}} - 1}. \quad (62)$$

Our strategy for tackling this problem is inspired by the proof technique in Corollary 6.7 of [3]. Let  $P$  be the data-dependent linear operator defined as,

$$[Ph](\mathbf{x}) \triangleq \frac{1}{K} \sum_{k=1}^K h(\mathbf{x}_k) g(\mathbf{x}, \mathbf{x}_k). \quad (63)$$

Let  $\mathcal{H}$  denote the Reproducing Kernel Hilbert Space associated with  $g$  and  $\langle \cdot, \cdot \rangle_{\mathcal{H}}$  be the dot product in  $\mathcal{H}$ . It is easy to verify that  $P$  is a **positive definite operator** in this space, i.e. it satisfies  $\langle h, Ph \rangle > 0$  for any  $h \in \mathcal{H}$  due to,

$$\langle h, Ph \rangle_{\mathcal{H}} = \langle h, \frac{1}{K} \sum_{k=1}^K h(\mathbf{x}_k) g(\cdot, \mathbf{x}_k) \rangle \quad (64)$$

$$= \frac{1}{K} \sum_{k=1}^K h(\mathbf{x}_k) \underbrace{\langle h, g(\cdot, \mathbf{x}_k) \rangle}_{h(\mathbf{x}_k)} \quad (65)$$

$$= \frac{1}{K} \sum_{k=1}^K h^2(\mathbf{x}_k) > 0, \quad (66)$$

where we used  $\langle h, g(\cdot, \mathbf{x}) \rangle = h(\mathbf{x})$  due to the reproducing property of  $\mathcal{H}$ . Since  $P$  is positive definite, there exist eigenfunctions  $\phi_j$  and eigenvalues  $\lambda_j \geq 0$  that satisfy  $[P\phi_j](\mathbf{x}) = \lambda_j \phi_j(\mathbf{x})$ . Plugging the definition of  $P$  into this identity yields,

$$\frac{1}{K} \sum_{k=1}^K \phi_j(\mathbf{x}_k) g(\mathbf{x}, \mathbf{x}_k) = \lambda_j \phi_j(\mathbf{x}). \quad (67)$$

In particular, evaluating the latter identity at the points  $\mathbf{x} \in \cup_{p=1}^K \{\mathbf{x}_p\}$  gives  $\frac{1}{K} \sum_{k=1}^K \phi_j(\mathbf{x}_k) g(\mathbf{x}_p, \mathbf{x}_k) = \lambda_j \phi_j(\mathbf{x}_p)$  for  $p = 1, \dots, K$ . Recalling that  $\mathbf{G}$  is evaluation of  $\frac{1}{K} g(\cdot, \cdot)$  at pairs of points across  $\cup_{k=1}^K \{\mathbf{x}_k\}$ , this identity can be expressed equivalently as,

$$\mathbf{G} \phi_j = \lambda_j \phi_j. \quad (68)$$

This implies  $\phi_j$  is an eigenvector of  $\mathbf{G}$  with corresponding eigenvalue of  $\lambda_j$  for any  $j$  that  $\mathbf{G} \phi_j \neq 0$ . Thus, by sorting  $\phi_j$  in non-increasing order of  $\lambda_j$ , and placing them for  $j = 1, \dots, K$  into the matrix  $\Phi$  and the diagonal matrix  $\Lambda$  respectively, we obtain,

$$\Phi = \mathbf{V}, \quad \Lambda = \mathbf{D}. \quad (69)$$

Since the eigenvectors of  $\mathbf{G}^t \Pi_{i=0}^t (\mathbf{G} + c_i \mathbf{I})^{-1}$  are the same as those of  $\mathbf{G}$  (adding a multiple of  $\mathbf{I}$  or applying matrix inversion do not change eigenvectors), and the eigenvectors of  $\mathbf{G}$  as showed in (68) are  $\Phi$ , we can write,

$$\mathbf{G}^t \Pi_{i=0}^t (\mathbf{G} + c_i \mathbf{I})^{-1} = \Phi^T \Lambda^t \Pi_{i=0}^t (\Lambda + c_i \mathbf{I})^{-1} \Phi. \quad (70)$$

On the other hand, using the same vector notation and recalling that  $\mathbf{g}$  is the evaluation of  $\frac{1}{K} g(\cdot, \mathbf{x}_k)$  at  $k = 1, \dots, K$ , we can express (67) as  $\phi_j^T \mathbf{g}_x = \lambda_j \phi_j(\mathbf{x})$ . Expressing this simultaneously for  $j = 1, \dots, K$  yields  $\Phi \mathbf{g}_x = \Lambda \phi_x$ , or equivalently

$$\mathbf{g}_x = \Phi^T \Lambda \phi_x, \quad (71)$$

where  $\phi_x \triangleq [\phi_1(\mathbf{x}), \dots, \phi_K(\mathbf{x})]$ . Plugging (70) and (71) with into (60) gives,

$$f_t^*(\mathbf{x}) = \mathbf{g}_x^T \mathbf{G}^t \Pi_{i=0}^t (\mathbf{G} + c_i \mathbf{I})^{-1} \mathbf{y}_0 \quad (72)$$

$$= \phi_x^T \Lambda \Phi \Phi^T \Lambda^t \Pi_{i=0}^t (\Lambda + c_i \mathbf{I})^{-1} \Phi \mathbf{y}_0 \quad (73)$$

$$= \phi_x^T \Lambda^{t+1} \Pi_{i=0}^t (\Lambda + c_i \mathbf{I})^{-1} \Phi \mathbf{y}_0. \quad (74)$$

Suppose  $g^\dagger$  is a positive definite kernel and let  $[P^\dagger h](x) \triangleq \frac{1}{K} \sum_{k=1}^K h(\mathbf{x}_k) g^\dagger(\mathbf{x}, \mathbf{x}_i)$ . We assume the operator  $P^\dagger$  shares the same eigenfunction as those of  $P$ , but varies in its eigenvalues  $\lambda_j^\dagger \geq 0$ , i.e.  $[P^\dagger \phi_j](\mathbf{x}) = \lambda_j^\dagger \phi_j(\mathbf{x})$ . Thus, by a similar argument, the solution of (61) can be written as,

$$f^\dagger(\mathbf{x}) = \phi_{\mathbf{x}}^T \mathbf{\Lambda}^\dagger (\mathbf{\Lambda}^\dagger + c_0 \mathbf{I})^{-1} \mathbf{\Phi} \mathbf{y}_0, \quad (75)$$

Thus in order to have  $f^\dagger = f_t^*$ , it is sufficient to have,

$$\mathbf{\Lambda}^{t+1} \Pi_{i=0}^t (\mathbf{\Lambda} + c_i \mathbf{I})^{-1} = \mathbf{\Lambda}^\dagger (\mathbf{\Lambda}^\dagger + c_0 \mathbf{I})^{-1}. \quad (76)$$

Since the matrices above are all diagonal, this can be expressed equivalently as,

$$\frac{\lambda_k^{t+1}}{\Pi_{i=0}^t (\lambda_k + c_i)} = \frac{\lambda_k^\dagger}{\lambda_k^\dagger + c_0}. \quad (77)$$

Solving in  $\lambda_k^\dagger$  yields,

$$\lambda_k^\dagger = c_0 \frac{1}{\frac{\Pi_{i=0}^t (\lambda_k + c_i)}{\lambda_k^{t+1}} - 1}. \quad (78)$$

Note that this is a valid solution for  $\lambda_k^\dagger$ , i.e. it satisfies the requirement  $\lambda_k^\dagger \geq 0$ . This is because  $\omega_k \triangleq \frac{\lambda_k^{t+1}}{\Pi_{i=0}^t (\lambda_k + c_i)}$  always satisfies<sup>6</sup>  $0 < \omega_k < 1$  and that the function  $\lambda_k^\dagger(\omega_k) \triangleq c_0 \frac{1}{\frac{1}{\omega_k} - 1}$  is well-defined ( $\omega_k \neq 0$ ) and is increasing when  $0 < \omega_k < 1$ .

---

<sup>6</sup>This is due to the conditions  $\lambda_k > 0$  (recall we assume  $\mathbf{G}$  is full-rank) and  $c_i > 0$ .

## C Proofs

**Proposition 1** *The variational problem (9) has a solution of the form,*

$$f^*(\mathbf{x}) = \mathbf{g}_{\mathbf{x}}^T (c\mathbf{I} + \mathbf{G})^{-1} \mathbf{y}. \quad (79)$$

See Appendix A for a proof.

**Proposition 2** *The following identity holds,*

$$\frac{1}{K} \sum_k (f^*(\mathbf{x}_k) - y_k)^2 = \frac{1}{K} \sum_k \left( z_k \frac{c}{c + d_k} \right)^2. \quad (80)$$

**Proof**

$$\frac{1}{K} (f^*(\mathbf{x}_k) - y_k)^2 \quad (81)$$

$$= \frac{1}{K} (\mathbf{g}_{\mathbf{x}_k}^T (c\mathbf{I} + \mathbf{G})^{-1} \mathbf{y} - y_k)^2 \quad (82)$$

$$= \frac{1}{K} \|\mathbf{G}(c\mathbf{I} + \mathbf{G})^{-1} \mathbf{y} - \mathbf{y}\|^2 \quad (83)$$

$$= \frac{1}{K} \|\mathbf{V}^T \mathbf{D}(c\mathbf{I} + \mathbf{D})^{-1} \mathbf{V} \mathbf{y} - \mathbf{y}\|^2, \quad (84)$$

which after exploiting rotation invariance property of  $\|\cdot\|$  and the fact that the matrix of eigenvectors  $\mathbf{V}$  is a rotation matrix, can be expressed as,

$$\frac{1}{K} (f^*(\mathbf{x}_k) - y_k)^2 \quad (85)$$

$$= \frac{1}{K} \|\mathbf{V}^T \mathbf{D}(c\mathbf{I} + \mathbf{D})^{-1} \mathbf{V} \mathbf{y} - \mathbf{y}\|^2 \quad (86)$$

$$= \frac{1}{K} \|\mathbf{V} \mathbf{V}^T \mathbf{D}(c\mathbf{I} + \mathbf{D})^{-1} \mathbf{V} \mathbf{y} - \mathbf{V} \mathbf{y}\|^2 \quad (87)$$

$$= \frac{1}{K} \|\mathbf{D}(c\mathbf{I} + \mathbf{D})^{-1} \mathbf{z} - \mathbf{z}\|^2 \quad (88)$$

$$= \frac{1}{K} \left\| (\mathbf{D}(c\mathbf{I} + \mathbf{D})^{-1} - \mathbf{I}) \mathbf{z} \right\|^2 \quad (89)$$

$$= \frac{1}{K} \sum_k \left( \frac{d_k}{c + d_k} - 1 \right)^2 z_k^2 \quad (90)$$

$$= \frac{1}{K} \sum_k \left( z_k \frac{c}{c + d_k} \right)^2, \quad (91)$$

□

**Proposition 3** *For any  $t \geq 0$ , if  $\|\mathbf{z}_i\| > \sqrt{K}\epsilon$  for  $i = 0, \dots, t$ , then,*

$$\|\mathbf{z}_t\| \geq a^t(\kappa) \|\mathbf{z}_0\| - \sqrt{K}\epsilon b(\kappa) \frac{a^t(\kappa) - 1}{a(\kappa) - 1}, \quad (92)$$

where,

$$a(x) \triangleq \frac{(r_0 - 1)^2 + x(2r_0 - 1)}{(r_0 - 1 + x)^2} \quad (93)$$

$$b(x) \triangleq \frac{r_0^2 x}{(r_0 - 1 + x)^2} \quad (94)$$

$$r_0 \triangleq \frac{1}{\sqrt{K}\epsilon} \|\mathbf{z}_0\|, \quad \kappa \triangleq \frac{d_{\max}}{d_{\min}}. \quad (95)$$

**Proof** We start from the identity we obtained in (17). By diving both sides of it by  $\sqrt{K}\epsilon$  we obtain,

$$\frac{1}{\sqrt{K}\epsilon} \mathbf{z}_t = \mathbf{D} \left( \frac{\alpha_t \sqrt{K}\epsilon}{\|\mathbf{z}_{t-1}\| - \sqrt{K}\epsilon} \mathbf{I} + \mathbf{D} \right)^{-1} \frac{1}{\sqrt{K}\epsilon} \mathbf{z}_{t-1}, \quad (96)$$

where,

$$d_{\min} \leq \alpha_t \leq d_{\max}. \quad (97)$$

Note that the matrix  $D(\frac{\alpha_t \sqrt{K\epsilon}}{\|z_{t-1}\| - \sqrt{K\epsilon}} I + D)^{-1}$  in the above identity is *diagonal* and its  $k$ 'th entry can be expressed as,

$$(D(\frac{\alpha_t \sqrt{K\epsilon}}{\|z_{t-1}\| - \sqrt{K\epsilon}} I + D)^{-1})[k, k] = \frac{d_k}{\frac{\alpha_t \sqrt{K\epsilon}}{\|z_{t-1}\| - \sqrt{K\epsilon}} + d_k} = \frac{1}{\frac{\frac{\alpha_t}{d_k}}{\frac{\|z_{t-1}\|}{\sqrt{K\epsilon}} - 1} + 1}. \quad (98)$$

Thus, as long as  $\|z_{t-1}\| > \sqrt{K\epsilon}$  we can get the following upper and lower bounds,

$$\frac{1}{\frac{\frac{d_{\max}}{d_{\min}}}{\frac{\|z_{t-1}\|}{\sqrt{K\epsilon}} - 1} + 1} \leq (D(\frac{\alpha_t \sqrt{K\epsilon}}{\|z_{t-1}\| - \sqrt{K\epsilon}} I + D)^{-1})[k, k] \leq \frac{1}{\frac{\frac{d_{\min}}{d_{\max}}}{\frac{\|z_{t-1}\|}{\sqrt{K\epsilon}} - 1} + 1}. \quad (99)$$

Putting the above fact beside recurrence relation of  $z_t$  in (96), we can bound  $\frac{1}{\sqrt{K\epsilon}} \|z_t\|$  as,

$$\frac{1}{\frac{\kappa}{r_{t-1}-1} + 1} r_{t-1} \leq r_t \leq \frac{1}{\frac{\frac{1}{\kappa}}{r_{t-1}-1} + 1} r_{t-1}, \quad (100)$$

where we used short hand notation,

$$\kappa \triangleq \frac{d_{\max}}{d_{\min}} \quad (101)$$

$$r_t \triangleq \frac{1}{\sqrt{K\epsilon}} \|z_t\|. \quad (102)$$

Note that  $\kappa$  is the *condition number* of the matrix  $G$  and by definition satisfies  $\kappa \geq 1$ . To further simplify the bounds, we use the inequality<sup>7</sup>,

$$\frac{1}{\frac{\frac{1}{\kappa}}{r_{t-1}-1} + 1} r_{t-1} \leq r_{t-1} \frac{(r_0 - 1)^2 + \frac{1}{\kappa}(2r_0 - 1)}{(r_0 - 1 + \frac{1}{\kappa})^2} - \frac{r_0^2 \frac{1}{\kappa}}{(r_0 - 1 + \frac{1}{\kappa})^2}, \quad (103)$$

and<sup>8</sup>,

$$\frac{1}{\frac{\kappa}{r_{t-1}-1} + 1} r_{t-1} \geq r_{t-1} \frac{(r_0 - 1)^2 + \kappa(2r_0 - 1)}{(r_0 - 1 + \kappa)^2} - \frac{r_0^2 \kappa}{(r_0 - 1 + \kappa)^2}. \quad (104)$$

For brevity, we introduce,

$$a(x) \triangleq \frac{(r_0 - 1)^2 + x(2r_0 - 1)}{(r_0 - 1 + x)^2} \quad (105)$$

$$b(x) \triangleq \frac{r_0^2 x}{(r_0 - 1 + x)^2}. \quad (106)$$

Therefore, the bounds can be expressed more concisely as,

$$a(\kappa) r_{t-1} - b(\kappa) \leq r_t \leq a(\frac{1}{\kappa}) r_{t-1} - b(\frac{1}{\kappa}). \quad (107)$$

Now since both  $r_{t-1} \triangleq \frac{1}{\sqrt{K\epsilon}} \|z_{t-1}\|$  and  $a(\kappa)$  or  $a(\frac{1}{\kappa})$  are non-negative, we can solve the recurrence<sup>9</sup> and obtain,

$$\frac{a^t(\kappa) r_0 - b(\kappa)}{a(\kappa) - 1} \leq r_t \leq \frac{a^t(\frac{1}{\kappa}) r_0 - b(\frac{1}{\kappa})}{a(\frac{1}{\kappa}) - 1}. \quad (108)$$

<sup>7</sup>This follows from concavity of  $\frac{x}{\frac{1}{\kappa} - x + 1}$  in  $x$  as long as  $x - 1 \geq 0$  (can be verified by observing that the second derivative of the function is negative when  $x - 1 \geq 0$  because  $\kappa > 1$  by definition). For any function  $f(x)$  that is concave on the interval  $[x, \bar{x}]$ , any line tangent to  $f$  forms an *upper* bound on  $f(x)$  over  $[x, \bar{x}]$ . In particular, we use the tangent at the end point  $\bar{x}$  to construct our bound. In our setting, this point which happens to be  $r_0$ . The latter is because  $r_t$  is a decreasing sequence (see beginning of Section 3.2) and thus its largest values is at  $t = 0$ .

<sup>8</sup>Similar to the earlier footnote, this follows from convexity of  $\frac{x}{\frac{\kappa}{x-1} + 1}$  in  $x$  as long as  $x - 1 \geq 0$  since  $\kappa > 1$  by definition. For any function  $f(x)$  that is convex on the interval  $[x, \bar{x}]$ , any line tangent to  $f$  forms an *lower* bound on  $f(x)$  over  $[x, \bar{x}]$ . In particular, we use the tangent at the end point  $\bar{x}$  to construct our bound, which as the earlier footnote, translate into  $r_0$ .

<sup>9</sup>More compactly, the problem can be stated as  $\alpha^\dagger r_{t-1} - b \leq r_t \leq \alpha r_{t-1} - b$ , where  $\alpha > 0$  and  $\alpha^\dagger > 0$ . Let's focus on  $r_t \leq \alpha r_{t-1} - b$ , as the other case follows by similar argument. Start from the base case  $r_1 \leq \alpha r_0 - b$ . Since  $\alpha > 0$ , we can multiply both sides by that and then add  $-b$  to both sides:  $\alpha r_1 - b \leq \alpha^2 r_0 - b(\alpha + 1)$ . On the other hand, looking at the recurrence  $r_t \leq \alpha r_{t-1} - b$  at  $t = 2$  yields  $r_2 \leq \alpha r_1 - b$ . Combining the two inequalities gives  $r_2 \leq \alpha^2 r_0 - b(\alpha + 1)$ . By repeating this argument we obtain the general case  $r_t \leq \alpha^t r_0 - b(\sum_{j=0}^{t-1} \alpha^j)$ .



□

**Proposition 4** Starting from  $\|\mathbf{y}_0\| > \sqrt{K}\epsilon$ , meaningful (non-collapsing solution) self-distillation is possible at least for  $\underline{t}$  rounds,

$$\underline{t} \triangleq \frac{\frac{\|\mathbf{y}_0\|}{\sqrt{K}\epsilon} - 1}{\kappa}. \quad (109)$$

**Proof** Recall that the assumption  $\|\mathbf{z}_t\| > \sqrt{K}\epsilon$  translates into  $r_t > 1$ . We now obtain a sufficient condition for  $r_t > 1$  by requiring a lower bound on  $r_t$  to be greater than one. For that purpose, we utilize the lower bound we established in (108),

$$\underline{r}_t \triangleq a^t(\kappa)r_0 - b(\kappa)\frac{a^t(\kappa) - 1}{a(\kappa) - 1}. \quad (110)$$

Setting the above to value 1 implies,

$$\underline{r}_t = 1 \Rightarrow t = \frac{\log\left(\frac{1-a(\kappa)+b(\kappa)}{b(\kappa)+r_0(1-a(\kappa))}\right)}{\log(a(\kappa))} = \frac{\log\left(\frac{1+\frac{\kappa-1}{r_0^2}}{1+\frac{\kappa-1}{r_0}}\right)}{\log\left(1 - \frac{(\frac{\kappa-1}{r_0} + \frac{1}{r_0})(\frac{\kappa-1}{r_0})}{(1+\frac{\kappa-1}{r_0})^2}\right)}. \quad (111)$$

Observe that,

$$\frac{\log\left(\frac{1+\frac{\kappa-1}{r_0^2}}{1+\frac{\kappa-1}{r_0}}\right)}{\log\left(1 - \frac{(\frac{\kappa-1}{r_0} + \frac{1}{r_0})(\frac{\kappa-1}{r_0})}{(1+\frac{\kappa-1}{r_0})^2}\right)} \geq \frac{r_0 - 1}{\kappa}, \quad (112)$$

Thus,

$$t \geq \frac{r_0 - 1}{\kappa} = \frac{\frac{\|\mathbf{z}_0\|}{\sqrt{K}\epsilon} - 1}{\kappa} = \frac{\frac{\|\mathbf{z}_0\|}{\sqrt{K}\epsilon} - 1}{\kappa} = \frac{\frac{\|\mathbf{y}_0\|}{\sqrt{K}\epsilon} - 1}{\kappa}. \quad (113)$$

□

**Theorem 5** Suppose  $\|\mathbf{y}_0\| > \sqrt{K}\epsilon$  and  $t \leq \frac{\|\mathbf{y}_0\|}{\kappa\sqrt{K}\epsilon} - \frac{1}{\kappa}$ . Then for any pair of diagonals of  $\mathbf{D}$ , namely  $d_j$  and  $d_k$ , with the condition that  $d_k > d_j$ , the following inequality holds.

$$\frac{\mathbf{B}_{t-1}[k, k]}{\mathbf{B}_{t-1}[j, j]} \geq \left( \frac{\frac{\|\mathbf{y}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_j}}{\frac{\|\mathbf{y}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_k}} \right)^t. \quad (114)$$

**Proof** We start with the definition of  $\mathbf{A}_t$  from (13) and proceed as,

$$\frac{\mathbf{A}_t[k, k]}{\mathbf{A}_t[j, j]} = \frac{1 + \frac{c_t}{d_j}}{1 + \frac{c_t}{d_k}}. \quad (115)$$

Since the derivative of the r.h.s. above w.r.t.  $c_t$  is non-negative as long as  $d_k \geq d_j$ , it is non-decreasing in  $c_t$ . Therefore, we can get a lower bound on r.h.s. using a lower bound on  $c_t$  (denoted by  $\underline{c}_t$ ),

$$\frac{\mathbf{A}_t[k, k]}{\mathbf{A}_t[j, j]} \geq \frac{1 + \frac{\underline{c}_t}{d_j}}{1 + \frac{\underline{c}_t}{d_k}}. \quad (116)$$

Also, since the assumption  $t \leq \frac{\|\mathbf{y}_0\|}{\kappa\sqrt{K}\epsilon} - \frac{1}{\kappa}$  guarantees non-collapse conditions  $c_t > 0$  and  $\|\mathbf{z}_t\| > \sqrt{K}\epsilon$ , we can apply (16) and have the following lower bound on  $c_t$

$$c_t \geq \frac{d_{\min}\sqrt{K}\epsilon}{\|\mathbf{z}_t\| - \sqrt{K}\epsilon}. \quad (117)$$

Since the r.h.s. (117) is decreasing in  $\|\mathbf{z}_t\|$ , the smallest value for the r.h.s. is attained by the largest value of  $\|\mathbf{z}_t\|$ . However, as  $\|\mathbf{z}_t\|$  is decreasing in  $t$  (see beginning of Section 3.2), its largest value is attained at  $t = 0$ . Putting these together we obtain,

$$c_t \geq \frac{d_{\min}\sqrt{K}\epsilon}{\|\mathbf{z}_0\| - \sqrt{K}\epsilon}. \quad (118)$$

Using the r.h.s. of the above as  $\underline{c}_t$  and applying it to (116) yields,

$$\frac{\mathbf{A}_t[k, k]}{\mathbf{A}_t[j, j]} \geq \frac{\frac{\|\mathbf{z}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_j}}{\frac{\|\mathbf{z}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_k}}. \quad (119)$$

Notice that both sides of the inequality are positive;  $\mathbf{A}_t$  based on its definition in (13) and r.h.s. by the fact that  $\|\mathbf{z}_0\| \geq \sqrt{K}\epsilon$ . Therefore, we can instantiate the above inequality at each distillation step  $i$ , for  $i = 0, \dots, t-1$ , and multiply them to obtain,

$$\prod_{i=0}^{t-1} \frac{\mathbf{A}_i[k, k]}{\mathbf{A}_i[j, j]} \geq \left( \frac{\frac{\|\mathbf{z}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_j}}{\frac{\|\mathbf{z}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_k}} \right)^t. \quad (120)$$

or equivalently,

$$\frac{\mathbf{B}_{t-1}[k, k]}{\mathbf{B}_{t-1}[j, j]} \geq \left( \frac{\frac{\|\mathbf{z}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_j}}{\frac{\|\mathbf{z}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_k}} \right)^t. \quad (121)$$

□

**Theorem 6** Suppose  $\|\mathbf{y}_0\| > \sqrt{K}\epsilon$ . Then the sparsity index  $S_{\mathbf{B}_{\underline{t}-1}}$  (where  $\underline{t} = \frac{\|\mathbf{y}_0\|}{\kappa\sqrt{K}\epsilon} - \frac{1}{\kappa}$  is number of guaranteed self-distillation steps before solution collapse) “decreases” in  $\epsilon$ , i.e. lower  $\epsilon$  yields higher sparsity.

Furthermore at the limit  $\epsilon \rightarrow 0$ , the sparsity index has the form,

$$\lim_{\epsilon \rightarrow 0} S_{\mathbf{B}_{\underline{t}-1}} = e^{\frac{d_{\min}}{\kappa} \min_{k \in \{1, 2, \dots, K-1\}} \left( \frac{1}{d_k} - \frac{1}{d_{k+1}} \right)}. \quad (122)$$

**Proof** We first show that the sparsity index is decreasing in  $\epsilon$ . We start from the definition of the sparsity index  $S_{\mathbf{B}_{\underline{t}-1}}$  in (24) which we repeat below,

$$S_{\mathbf{B}_{\underline{t}-1}} = \min_{k \in \{1, 2, \dots, K-1\}} \left( \frac{\frac{\|\mathbf{y}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_k}}{\frac{\|\mathbf{y}_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_{k+1}}} \right)^{\frac{\|\mathbf{y}_0\|}{\kappa\sqrt{K}\epsilon} - \frac{1}{\kappa}}. \quad (123)$$

For brevity, we define base and exponent as,

$$b \triangleq \frac{m + \frac{d_{\min}}{d_k}}{m + \frac{d_{\min}}{d_{k+1}}} \quad (124)$$

$$p \triangleq \frac{m}{\kappa} \quad (125)$$

$$m \triangleq \frac{\|\mathbf{y}_0\|}{\sqrt{K}\epsilon} - 1, \quad (126)$$

so that,

$$S_{\mathbf{B}_{\underline{t}-1}}(\epsilon) = b^p. \quad (127)$$

The derivative is thus,

$$\frac{d}{d\epsilon} S_{\mathbf{B}_{\underline{t}-1}} \quad (128)$$

$$= \frac{d S_{\mathbf{B}_{\underline{t}-1}}}{dm} \frac{dm}{d\epsilon} \quad (129)$$

$$= \left( b^p \left( \frac{p b_m}{b} + p_m \log(b) \right) \right) \left( \frac{dm}{d\epsilon} \right) \quad (130)$$

$$= b^p \left( \frac{p b_m}{b} + p_m \log(b) \right) \left( -\frac{1}{2\epsilon} (m+1) \right) \quad (131)$$

$$= b^p \left( \frac{p}{m + \frac{d_{\min}}{d_k}} - \frac{p}{m + \frac{d_{\min}}{d_{k+1}}} + \frac{1}{\kappa} \log(b) \right) \left( -\frac{1}{2\epsilon} (m+1) \right) \quad (132)$$

$$= \frac{b^p}{\kappa} \left( \frac{m}{m + \frac{d_{\min}}{d_k}} - \frac{m}{m + \frac{d_{\min}}{d_{k+1}}} + \log(b) \right) \left( -\frac{1}{2\epsilon} (m+1) \right) \quad (133)$$

$$= \frac{b^p}{\kappa} \left( \frac{1}{1 + \frac{d_{\min}}{m d_k}} - \frac{1}{1 + \frac{d_{\min}}{m d_{k+1}}} + \log(b) \right) \left( -\frac{1}{2\epsilon} (m+1) \right) \quad (134)$$

$$= \frac{b^p}{\kappa} \left( \frac{1}{1 + \frac{d_{\min}}{m d_k}} - \frac{1}{1 + \frac{d_{\min}}{m d_{k+1}}} + \log\left(\frac{1 + \frac{d_{\min}}{m d_k}}{1 + \frac{d_{\min}}{m d_{k+1}}}\right) \right) \left( -\frac{1}{2\epsilon} (m+1) \right) \quad (135)$$

$$= \frac{b^p}{\kappa} \left( \frac{1}{1 + \frac{d_{\min}}{m d_k}} + \log\left(1 + \frac{d_{\min}}{m d_k}\right) - \frac{1}{1 + \frac{d_{\min}}{m d_{k+1}}} - \log\left(1 + \frac{d_{\min}}{m d_{k+1}}\right) \right) \left( -\frac{1}{2\epsilon} (m+1) \right). \quad (136)$$

We now focus on the first parentheses. Define the function  $e(x) \triangleq \frac{1}{x} + \log(x)$ . Thus we can write the contents in the first parentheses more compactly,

$$\frac{1}{1 + \frac{d_{\min}}{m d_k}} + \log\left(1 + \frac{d_{\min}}{m d_k}\right) - \frac{1}{1 + \frac{d_{\min}}{m d_{k+1}}} - \log\left(1 + \frac{d_{\min}}{m d_{k+1}}\right) \quad (137)$$

$$= e\left(1 + \frac{d_{\min}}{m d_k}\right) - e\left(1 + \frac{d_{\min}}{m d_{k+1}}\right). \quad (138)$$

However,  $e'(x) = \frac{x-1}{x^2}$ , thus when  $x > 1$  the function  $e'(x)$  is positive. Consequently, when  $x > 1$   $e(x)$  is increasing. In fact, since both  $\frac{d_{\min}}{m d_k}$  and  $\frac{d_{\min}}{m d_{k+1}}$  are positive, the arguments of  $e$  satisfy the condition of being greater than 1 and thus  $e$  is increasing. On the other hand, since  $d_{k+1} > d_k$  it follows that  $1 + \frac{d_{\min}}{m d_k} > 1 + \frac{d_{\min}}{m d_{k+1}}$ , and thus by leveraging the fact that  $e$  is increasing we obtain  $e\left(1 + \frac{d_{\min}}{m d_k}\right) > e\left(1 + \frac{d_{\min}}{m d_{k+1}}\right)$ . Finally by plugging the definition of  $e$  we obtain,

$$\frac{1}{1 + \frac{d_{\min}}{m d_k}} + \log\left(1 + \frac{d_{\min}}{m d_k}\right) > \frac{1}{1 + \frac{d_{\min}}{m d_{k+1}}} + \log\left(1 + \frac{d_{\min}}{m d_{k+1}}\right). \quad (139)$$

It is now easy to determine the sign of  $\frac{d}{d\epsilon} S$  as shown below,

$$\begin{aligned} & \frac{d}{d\epsilon} S_{B_{t-1}} \quad (140) \\ = & \underbrace{\frac{b^p}{\kappa}}_{\text{positive}} \underbrace{\left( \frac{1}{1 + \frac{d_{\min}}{m d_k}} + \log\left(1 + \frac{d_{\min}}{m d_k}\right) - \frac{1}{1 + \frac{d_{\min}}{m d_{k+1}}} - \log\left(1 + \frac{d_{\min}}{m d_{k+1}}\right) \right)}_{\text{positive}} \underbrace{\left( -\frac{1}{2\epsilon}(m+1) \right)}_{\text{negative}} \quad (141) \end{aligned}$$

By showing that  $\frac{d}{d\epsilon} S_{B_{t-1}} < 0$  we just proved  $S_{B_{t-1}}$  is decreasing in  $\epsilon$ .

We now focus on the limit case  $\epsilon \rightarrow 0$ . First note due to the identity  $m = \frac{\|y_0\|}{\sqrt{K}\epsilon} - 1$  we have the following identity,

$$\lim_{\epsilon \rightarrow 0} \min_{k \in \{1, 2, \dots, K-1\}} \left( \frac{\frac{\|y_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_k}}{\frac{\|y_0\|}{\sqrt{K}\epsilon} - 1 + \frac{d_{\min}}{d_{k+1}}} \right)^{\frac{\|y_0\|}{\sqrt{K}\epsilon} - \frac{1}{\kappa}} \quad (142)$$

$$= \lim_{m \rightarrow \infty} \min_{k \in \{1, 2, \dots, K-1\}} \left( \frac{m + \frac{d_{\min}}{d_k}}{m + \frac{d_{\min}}{d_{k+1}}} \right)^{\frac{1}{\kappa} m}. \quad (143)$$

Further, since pointwise minimum of continuous functions is also a continuous function, we can move the limit inside the minimum,

$$\lim_{m \rightarrow \infty} \min_{k \in \{1, 2, \dots, K-1\}} \left( \frac{m + \frac{d_{\min}}{d_k}}{m + \frac{d_{\min}}{d_{k+1}}} \right)^{\frac{1}{\kappa} m} \quad (144)$$

$$= \min_{k \in \{1, 2, \dots, K-1\}} \lim_{m \rightarrow \infty} \left( \frac{m + \frac{d_{\min}}{d_k}}{m + \frac{d_{\min}}{d_{k+1}}} \right)^{\frac{1}{\kappa} m} \quad (145)$$

$$= \min_{k \in \{1, 2, \dots, K-1\}} e^{\frac{d_{\min}}{d_k} - \frac{d_{\min}}{d_{k+1}} \frac{1}{\kappa}} \quad (146)$$

$$= \min_{k \in \{1, 2, \dots, K-1\}} e^{\frac{d_{\min}}{\kappa} \left( \frac{1}{d_k} - \frac{1}{d_{k+1}} \right)} \quad (147)$$

$$= e^{\frac{d_{\min}}{\kappa} \min_{k \in \{1, 2, \dots, K-1\}} \left( \frac{1}{d_k} - \frac{1}{d_{k+1}} \right)}, \quad (148)$$

where in (146) we used the identity  $\lim_{x \rightarrow \infty} f(x)^{g(x)} = e^{\lim_{x \rightarrow \infty} (f(x)-1)(g(x))}$  and in (148) we used the fact that  $e^{\frac{d_{\min}}{\kappa} x}$  is monotonically increasing in  $x$  (because  $\frac{d_{\min}}{\kappa} > 0$ ).

□

## D More on Experiments

### D.1 Setup Details

We used Adam optimizer with learning rates of 0.001 and 0.0001 for CIFAR-10 and CIFAR-100, respectively. They are trained up to 64000 steps with batch size equal to 16 and 64 for CIFAR-10 and CIFAR-100, respectively. In all the experiments, we slightly regularize the training by weight decay regularization added to the fitting loss with its coefficient set to 0.0001 and 0.00005 for CIFAR-10 and CIFAR-100, respectively. Training and test is performed on the standard (50000 train-10000 test) split of the CIFAR dataset. Most of the experiments are conducted using Resnet-50 [12] and CIFAR-10 and CIFAR-100 datasets [18]. However, we briefly validate our results on VGG-16 [30] too.

### D.2 $\ell_2$ Loss on Neural Network Predictions

Figure 4 shows the full results on CIFAR-10 and Resnet-50. The train and test accuracies have already been discussed in the main paper and are copied here to facilitate comparison. However, in this subsection, we demonstrated the loss of the trained model at all steps with respect to the original ground truth data too. This may help establish an intuition on how self-distillation is regularizing the training on the original data. Looking at the train loss we can see it first drops as the regularization is amplified and then increases while the model under-fits. This, again, suggests that the mechanism that self-distillation employs for regularization is different from early stopping. For CIFAR-100 the results in Figure 5 show a similar trend.

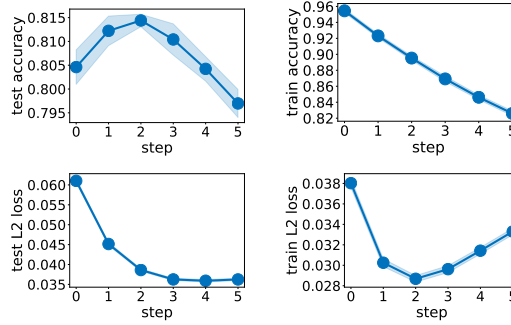


Figure 4: Self distillation results with  $\ell_2$  loss of neural network predictions for Resnet-50 and CIFAR-10

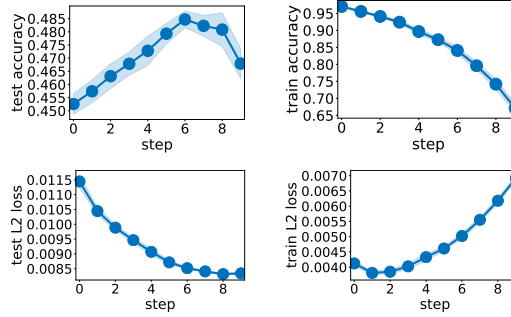


Figure 5: Self distillation results with  $\ell_2$  loss of neural network predictions for Resnet-50 and CIFAR-100

### D.3 Self-distillation on Hard Labels

One might wonder how self-distillation would perform if we replace the neural network (soft) predictions with hard labels. In other words, the teacher’s predictions are turned into one-hot-vector via `argmax` and they are treated like a dataset with augmented labels. Of-course, since the model is already over-parameterized and trained close to interpolation regime only a small fraction of labels will change. Figures 6 and 7 show the results of self distillation using cross entropy loss on labels predicted by the teacher model. Surprisingly, self-distillation improves the performance here too. This observation may be related to learning under noisy dataset and calls for more future work on this interesting case.

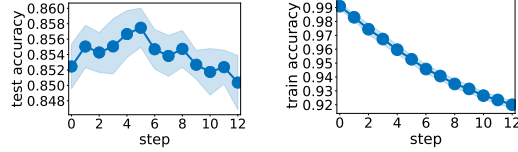


Figure 6: Self distillation results with cross entropy loss on hard labels for Resnet-50 and CIFAR-10

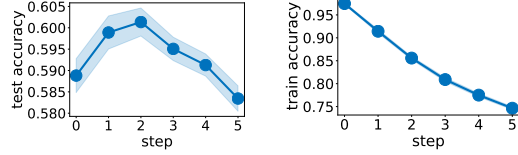


Figure 7: Self distillation results with cross entropy loss on hard labels for Resnet-50 and CIFAR-100

#### D.4 Self-Distillation versus Early Stopping.

By looking at the fall of the training accuracy over self-distillation round, one may wonder if early stopping (in the sense of choosing a larger error tolerance  $\epsilon$  for training) would lead to similar test performance. However, in Section 3.4 we discussed that self-distillation and early stopping have different regularization effects. Here we try to verify that. Specifically, we record the training loss value at the end of each self-distillation round. We then train a batch of models from scratch until each batch converges to one of the recorded loss values. If the regularization induced by early stopping was the same as self-distillation, then we should have seen similar test performance between a self-distilled model that achieves a specific loss value on the original training labels, and a model that stops training as soon as it reaches the same level of error. However, Figure 8 verifies that these two have different regularization effects.

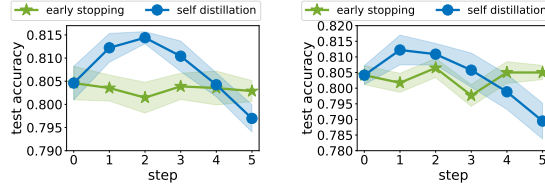


Figure 8: Self-distillation compared to early stopping for Resnet50 and CIFAR-10 using  $\ell_2$  and cross entropy loss, respectively.

**Self-Distillation on Other Networks.** Figure 9 shows the performance of  $\ell_2$  distillation on CIFAR-100 using VGG network. This result aims to show that the theory and empirical findings are not dependent to a specific structure and apply to architectures beyond Resnet.

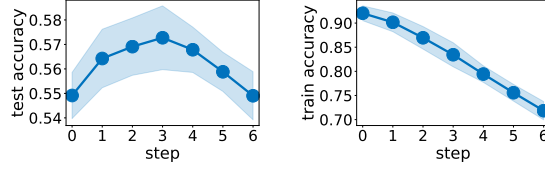


Figure 9: Self-distillation with  $\ell_2$  loss using VGG16 Network on CIFAR-100.

## E Mathematica Code To Reproduce Illustrative Example

```

x = (Table[i, {i, -5, 5}]/5 + 1)/2;
y = Sin[x*2*Pi] +
  RandomVariate[NormalDistribution[0, 0.5], Length[x]]
ListPlot[y]

(* UNCOMMENT IF YOU WISH TO USE EXACT SAME RANDOM SAMPLES IN THE PAPER *)
(* y = {0.38476636465198066',
  1.2333967683416893', 1.33232242218057',
  0.6920159488889518', -0.29756145531871736', -0.24189291901377769', \
-0.7964485769175675', -0.9616480167034174', -0.49672509509916934', \
-0.3469066003991437', 0.5589512650600734'}; *)

(***** PLOT GREEN'S FUNCTION g0(X,T) FOR OPERATOR d^4/dx^4 *****)

g0 = 1/6*Max[{(T - X)^3, 0}] - 1/6*T*(1 - X)*(T^2 - 2*X + X^2);
ContourPlot[g0, {X, 0, 1}, {T, 0, 1}]
Plot3D[g0, {X, 0, 1}, {T, 0, 1}]

(**** COMPUTE g AND G *****)

G = Table[
  g0 /. X -> ((i/5 + 1)/2) /. T -> ((j/5 + 1)/2), {i, -5, 5}, {j, -5,
  5}];
g = Transpose[{Table[g0 /. T -> ((j/5 + 1)/2), {j, -5, 5}]}];

(**** PLOT GROUND-TRUTH FUNCTION (ORANGE) AND OVERFIT FUNCTION \
(BLUE) *****)
FNoReg = (Transpose[g].Inverse[
  G + 0.0000000001*IdentityMatrix[Length[x]]].Transpose[{y}])[1,
  1];
pts = Table[{x[[i]], y[[i]]}, {i, 1, Length[x]}];
Show[{ListPlot[pts], Plot[{FNoReg, Sin[X*2*Pi]}, {X, 0, 1}]}]

(**** PARAMETERS *****)
MaxIter = 10;
eps = 0.045;

(**** SUBROUTINES *****)
Loss[G_, yin_, c_] := Module[
  {t = (G.Inverse[c*IdentityMatrix[Length[yin]] + G] -
  IdentityMatrix[Length[x]]).yin},
  Total[Flatten[t]^2]/Length[yin]
];

FindRootsC[f_, c_] := Module[
  {Sol = Quiet[Solve[f == 0, c]], Sel},
  Sel = Select[
    c /. Sol, (Abs[Im[#]] < 0.00000001) && # > 0.00000001 &]
];

```

```

];

(***** MAIN *****)

(* Initialization *)
y0 = Transpose[{y}];
ycur = y0;
B = IdentityMatrix[Length[x]];
FunctionSequence = {};
ASequence = {};
BSequence = {};

(* Self-Distillation Loop *)
For[i = 1; i < MaxIter, i++,
  Print["Iteration ", i];
  Print["Norm[y]=", Norm[ycur]];
  L = Loss[G, ycur, c];
  RootsC = FindRootsC[L - eps, c];
  Switch [Length[RootsC], 0, (Print["No Root"]; Break[]), 1,
    Print["Found Unique Root c=", RootsC[[1]] ]];
  (* Now that root is unique *)
  RootC = RootsC[[1]];
  Print["Achieved Loss Value ", Loss[G, ycur, RootC]];
  U = G.Inverse[G + RootC*IdentityMatrix[Length[ycur]]];
  A = DiagonalMatrix[Eigenvalues[U]];
  f = (Transpose[g].Inverse[
    G + RootC*IdentityMatrix[Length[ycur]]].ycur)[[1, 1]];
  B = B.A;
  ycur = U.ycur;

  FunctionSequence = Append[FunctionSequence, f];
  ASequence = Append[ASequence, Diagonal[A]];
  BSequence = Append[BSequence, Diagonal[B]];
]

If[i == MaxIter, Print["Max Iterations Reached!"]]

Plot[FunctionSequence, {X, 0, 1}]
BarChart[ASequence, ChartStyle -> "DarkRainbow", AspectRatio -> 0.2,
  ImageSize -> Full]
BarChart[BSequence, ChartStyle -> "DarkRainbow", AspectRatio -> 0.2,
  ImageSize -> Full]

```

## F Python Implementation

Implementing self-distillation is quite straight forward provided with merely a customized loss that replaces the ground-truth labels with teacher predictions. Here, we provide a Tensorflow implementation of the self-distillation loss function:

```
1 def self_distillation_loss(labels, logits, model, reg_coef,
2                             teacher=None, data=None):
3     if teacher is None:
4         main_loss = tf.reduce_mean(tf.squared_difference(
5             labels, tf.nn.softmax(logits)))
6     else:
7         main_loss = tf.reduce_mean(tf.squared_difference(
8             tf.nn.softmax(teacher(data)), tf.nn.softmax(logits)))
9     reg_loss = reg_coef*tf.add_n([tf.nn.l2_loss(w) for w
10                                in model.trainable_weights])
11     total_loss = main_loss + reg_loss
12     return total_loss
```

The following snippet also demonstrates how one can use the above loss function to train a neural network using self-distillation.

```
1 def self_distillation_train(model, train_dataset, optimizer,
2                             reg_coef, epochs, teacher=None):
3     for epoch in range(epochs):
4         for iter, (x_batch_train,
5                    y_batch_train) in enumerate(train_dataset):
6             with tf.GradientTape() as tape:
7                 logits = model(x_batch_train, training=True)
8                 loss_value =
9                     self_distillation_loss(y_batch_train, logits, model,
10                                           reg_coef, teacher, x_batch_train)
11                 grads = tape.gradient(loss_value, model.trainable_weights)
12                 optimizer.apply_gradients(
13                     zip(grads, model.trainable_weights))
14     return model
15
16 teacher = None
17 reg_coef=1e-4
18 epochs=30
19 for step in range(distillation_steps):
20     model = get_resnet_model()
21     optimizer = keras.optimizers.Adam(learning_rate=learning_rate)
22     model = self_distillation_train(
23         model, train_dataset, optimizer, reg_coef, epochs, teacher)
24     teacher = model
```