

1 We thank all reviewers. It was very encouraging to read your opinion on the *strengths* of this work: importance of the
2 problem, lack of prior work, novelty of the analysis, and consistency between our theory with empirical observation.

3 **[R2] Positive about this work. Minor typos.** We are very happy that you liked our contributions and found our
4 results intriguing. We will surely fix typos in the final version.

5 **[R3] In practice people use both predictions and original labels, but this paper only uses predictions.** Thank
6 for pointing out. In the literature, “self-distillation” refers to a host of related ideas. We have adopted the variant
7 proposed by the well-cited work of [1], which only uses predictions. Incidentally, [1] compared training with pseudo
8 labels generated merely from teacher’s predictions versus its combination with the ground truth labels (respectively BAN
9 and BAN+L in their Table 2), and observed that BAN outperforms BAN+L. What you are referring to is closer to BAN+L
10 (unfortunately you did not specify a reference so that we could be more precise here). That being said, as self-distillation
11 is a new area, there is not yet enough evidence to claim one variant is better than others. As such, we do not think the
12 variant we studied is less qualified than others, if not more (due to BAN+L observation of [1]).

13 Regarding your interesting question of whether blending the original labels into the pseudo-labels can prevent collapse, it
14 indeed does that, but at the cost of undoing some of the regularization benefits of self-distillation (consistent with BAN+L
15 vs BAN observation of [1]). The emphasis on the original labels facilitates overfitting to those labels and diminishes the
16 regularization effect. Note that the collapse is *not a practical concern*, as one should stop even earlier than that when
17 over-regularization begins. The latter can be detected by measuring test performance on the validation set.

18 [1] T. Furlanello et al. “Born-Again Neural Networks”, *ICML 2018*.

19 **[R3] What is the purpose of lower bound on the number of rounds of self distillation.** To be clear, this is the
20 lower bound on the number of self-distillation rounds before the solution *collapses* to zero. This quantity is critical
21 in determining the ultimate strength of regularization that self-distillation can achieve. The reason is, as proved in
22 Sections 3.3 and 3.5, the sparsity level enhances with each self-distillation round (as long as collapse is not reached)
23 [line 188-194]. Once the solution collapses, nothing interesting happens from that point on [line 154-156]. Thus, the
24 highest achievable sparsity is right before the collapse, which is determined by the lower bound you mentioned.

25 **[R3] What is the purpose of bounding S_B in Section 3.5?** The goal of the paper is to understand the implicit
26 regularization of self-distillation. Our analysis reveals that this regularization leads to a sparser matrix B . S_B is simply
27 a way to quantify such sparsity; it takes the diagonal matrix B and returns a single number. Since we showed the ratio
28 between pairs of diagonals of B change monotonically in t [line 188-191], one can see that the sparsity index S_B
29 increases in t , hence showing self-distillation enhancing the sparsity. We will clarify this in the revision.

30 In addition, S_B is used to analyze how ϵ affects the achievable sparsity level. Theorem 6 reveals that being near the
31 interpolation regime can enhance the regularization effect of self-distillation (i.e. leading to sparser representation).

32 **[R1] In practice it is unlikely that anybody would do self-distillation, as it is easier to just set the regularization
33 parameter better in the first place.** There seems to be a *misunderstanding* here, which has caused the result to
34 seem trivial. While it is true that one can “just set the regularization parameter better in the first place”, our results
35 show that it is *not possible* to achieve the regularization effect due to self-distillation in this way. Increasing the ridge
36 regularization coefficient c will scale all eigenvalues of the kernel by the same factor. Such uniform scaling does not
37 change the sparsity pattern of the eigenvalues. In contrast, self-distillation exponentiates the eigenvalues by the number
38 of steps t , which results in a non-uniform scaling. This allows to shrink some directions more than others. We add
39 that increasing c without any self-distillation is similar to regularization by early stopping. We have discussed this in
40 Section 3.4 and emphasized that these two regularization schemes behave very differently.

41 **[R1] Is theoretical contribution of RKHS analysis alone interesting?** Our results give insight into the behavior of
42 self-distillation in the RKHS setting, which on its own is an interesting family of learning problems. The dynamics of
43 self-distillation in RKHS setting follows a nonlinear recurrence for which there is no closed form solution, and this
44 *significantly* complicates the analysis. We have been able to fully characterize the evolution of self-distillation in the
45 RKHS setting and prove how the spectrum of the kernel evolves in an intriguing way that leads to sparsity. The fact that
46 self-distillation promotes sparsity has *never* been noticed before, let alone proving it. We will clarify this in the revision.

47 **[R1] There is a gap between theory (RKHS/NTK) and practical neural networks.** We agree that a mathematically
48 rigorous characterization of self-distillation for deep neural networks - comparable to our results in the RKHS setting
49 - would be an exciting contribution. However, please note that as neural networks are highly nonlinear, almost *all*
50 theoretical developments *require* some simplifying assumptions to keep the analysis tractable.