

## Supplementary Material

We define here some notation in addition to that of Section 3 in the main text. We denote by  $\ell_i$  the per-instance loss,

$$L^1(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^\top \mathbf{x}_i), \quad (35)$$

$$\ell_i(u) = -y_i \log \sigma(u) - (1 - y_i) \log(1 - \sigma(u)) - \ell_i^*, \quad (36)$$

where  $\ell_i^*$  are constants chosen such that the minimum of  $\ell_i$  is 0, namely  $\ell_i^* = -y_i \log y_i - (1 - y_i) \log(1 - y_i)$ .

Slightly abusing notation, we write  $L(\tau) = L^1(\mathbf{w}(\tau)) = L(\mathbf{W}_1(\tau), \dots, \mathbf{W}_N(\tau))$  for the objective value at time  $\tau$ .

Finally, for a full-rank matrix  $\mathbf{A} \in \mathbb{R}^{d \times m}$  ( $m \geq 1$ ), we denote by  $\mathbf{P}_\mathbf{A} \in \mathbb{R}^{d \times d}$  the matrix of projection onto the span of  $\mathbf{A}$ ,

$$\mathbf{P}_\mathbf{A} = \begin{cases} \mathbf{I}, & m \geq d, \\ \mathbf{A}(\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top, & m < d. \end{cases} \quad (37)$$

### A. Properties of the Cross-Entropy Loss

**Theorem A.1** (Gradient). *The gradient of the cross-entropy loss (35) takes the form*

$$\nabla L^1(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n (\sigma(\mathbf{w}^\top \mathbf{x}_i) - y_i) \cdot \mathbf{x}_i. \quad (38)$$

*It always lies in the data span,  $\nabla L^1(\mathbf{w}) \in \text{span}(\mathbf{X})$ .*

*Proof.* Straightforward calculation.  $\square$

**Theorem A.2** (Global minima). *The global minimum of the cross-entropy loss (35) is 0 and the set of global minimisers is*

$$\{\mathbf{w} \in \mathbb{R}^d : \mathbf{X}^\top \mathbf{w} = \mathbf{X}^\top \mathbf{w}_*\}. \quad (39)$$

*Proof.* We know that  $L^1 \geq 0$  and  $L^1(\mathbf{w}_*) = 0$ , so 0 is the optimal objective value, and the set of global optima consists of all  $\mathbf{w}$  such that  $L^1(\mathbf{w}) = 0$ . The last condition is equivalent to  $\forall_i : \ell_i(\mathbf{w}) = 0$ , which in turn is equivalent to  $\forall_i : \sigma(\mathbf{w}^\top \mathbf{x}_i) = \sigma(\mathbf{w}_*^\top \mathbf{x}_i)$ . By monotonicity of  $\sigma$ , this is further equivalent to  $\forall_i : \mathbf{w}^\top \mathbf{x}_i = \mathbf{w}_*^\top \mathbf{x}_i$ , which is a restatement of (39).  $\square$

**Theorem A.3** (Restricted strong convexity). *Assume  $\mathbf{X}$  is full-rank. For any sublevel set  $\mathcal{W} = \{\mathbf{w} : L^1(\mathbf{w}) \leq l\}$ ,*

*there exists  $\mu > 0$  such that*

$$L^1(\mathbf{v}) \geq L^1(\mathbf{w}) + \nabla L^1(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) + \frac{\mu}{2} \|\mathbf{v} - \mathbf{w}\|^2 \quad (40)$$

*for all  $\mathbf{w}, \mathbf{v} \in \mathcal{W}$  such that  $\mathbf{v} - \mathbf{w} \in \text{span}(\mathbf{X})$ .*

*Proof.* Consider the 2nd-order Taylor expansion of  $L^1$  around  $\mathbf{w}$ ,

$$L^1(\mathbf{v}) = L^1(\mathbf{w}) + \nabla L^1(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) + \frac{1}{2} (\mathbf{v} - \mathbf{w})^\top [\nabla^2 L^1(\bar{\mathbf{w}})] (\mathbf{v} - \mathbf{w}), \quad (41)$$

where  $\nabla^2 L^1(\bar{\mathbf{w}})$  is the Hessian of  $L^1$  evaluated at  $\bar{\mathbf{w}}$ , a point lying between  $\mathbf{v}$  and  $\mathbf{w}$ . A straightforward calculation shows that the Hessian takes the form

$$\nabla^2 L^1(\bar{\mathbf{w}}) = \mathbf{X} \mathbf{D}_{\bar{\mathbf{w}}} \mathbf{X}^\top, \quad (42)$$

where

$$\mathbf{D}_{\bar{\mathbf{w}}} = \text{diag}[\sigma(\bar{\mathbf{w}}^\top \mathbf{x}_1)(1 - \sigma(\bar{\mathbf{w}}^\top \mathbf{x}_1)), \dots, \sigma(\bar{\mathbf{w}}^\top \mathbf{x}_n)(1 - \sigma(\bar{\mathbf{w}}^\top \mathbf{x}_n))]. \quad (43)$$

We will now show that there is a constant  $\omega > 0$  such that

$$\sigma(\bar{\mathbf{w}}^\top \mathbf{x}_i)(1 - \sigma(\bar{\mathbf{w}}^\top \mathbf{x}_i)) \geq \omega \quad (44)$$

for all  $\bar{\mathbf{w}} \in \mathcal{W}$  and  $i \in \{1, \dots, n\}$ , so that we can claim  $\mathbf{D}_{\bar{\mathbf{w}}} \succeq \omega \mathbf{I}$ , or consequently  $\nabla^2 L^1(\bar{\mathbf{w}}) \succeq \omega \mathbf{X} \mathbf{X}^\top$ .

Let  $\mathbf{w} \in \mathcal{W}$ . The bound on  $L^1(\mathbf{w})$  implies a bound on  $\ell_i(\mathbf{w}^\top \mathbf{x}_i)$  for all  $i$ ,

$$\ell_i(\mathbf{w}^\top \mathbf{x}_i) \leq n L^1(\mathbf{w}) \leq nl. \quad (45)$$

Because  $\ell_i$  is convex and  $\ell_i(u) \rightarrow \infty$  as  $u \rightarrow \pm\infty$ , we know that  $\ell_i^{-1}((-\infty, nl])$  is a bounded interval, and the finite union  $\cup_{i=1}^n \ell_i^{-1}((-\infty, nl])$  is also a bounded interval, whose size depends only on  $nl$  and the data. Hence, there exists  $K > 0$  such that  $\mathbf{w}^\top \mathbf{x}_i \in [-K, K]$  for all  $\mathbf{w} \in \mathcal{W}$  and  $i \in \{1, \dots, n\}$ . The existence of  $\omega > 0$  satisfying (44) follows.

Now, let us apply  $\nabla^2 L^1(\bar{\mathbf{w}}) \succeq \omega \mathbf{X} \mathbf{X}^\top$  to lower-bound (41):

$$L^1(\mathbf{v}) \geq L^1(\mathbf{w}) + \nabla L^1(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) + \frac{\omega}{2} (\mathbf{v} - \mathbf{w})^\top \mathbf{X} \mathbf{X}^\top (\mathbf{v} - \mathbf{w}). \quad (46)$$

Consider two cases. If  $n \geq d$ ,  $\mathbf{X}\mathbf{X}^\top$  is full-rank and  $\mathbf{X}\mathbf{X}^\top \geq \lambda_{\min}\mathbf{I}$  holds, where  $\lambda_{\min} > 0$  is the smallest eigenvalue of  $\mathbf{X}\mathbf{X}^\top$ . Combined with (46), this proves the claim for  $n \geq d$  and  $\mu = \omega\lambda_{\min}$ .

If  $n < d$ ,  $\mathbf{X}^\top\mathbf{X}$  is full rank. We can use the assumption  $\mathbf{v} - \mathbf{w} \in \text{span}(\mathbf{X})$  to deduce

$$\begin{aligned} \|\mathbf{v} - \mathbf{w}\|^2 &= \|\mathbf{P}_{\mathbf{X}}(\mathbf{v} - \mathbf{w})\|^2 \\ &= (\mathbf{v} - \mathbf{w})^\top \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1} \mathbf{X}^\top(\mathbf{v} - \mathbf{w}) \quad (47) \\ &\leq \lambda_{\max}(\mathbf{v} - \mathbf{w})^\top \mathbf{X}\mathbf{X}^\top(\mathbf{v} - \mathbf{w}), \end{aligned}$$

where  $\lambda_{\max} > 0$  is the largest eigenvalue of  $(\mathbf{X}^\top\mathbf{X})^{-1}$ . Combined with (46), this proves the claim for  $n < d$  and  $\mu = \omega/\lambda_{\max}$ .  $\square$

**Corollary A.1** (Restricted Polyak-Lojasiewicz). *Assume  $\mathbf{X}$  is full-rank. For any sublevel set  $\mathcal{W} = \{\mathbf{w} : L^1(\mathbf{w}) \leq l\}$ , there exists  $c > 0$  such that*

$$cL^1(\mathbf{w}) \leq \frac{1}{2} \|\nabla L^1(\mathbf{w})\|^2 \quad (48)$$

for all  $\mathbf{w} \in \mathcal{W}$ .

*Proof.* Let  $\mathbf{w} \in \mathcal{W}$ . (If  $\mathcal{W}$  is empty, the claim is trivially true.) Theorem A.3 applied to  $\mathcal{W}$  implies that for some  $\mu > 0$ ,

$$L^1(\mathbf{v}) \geq L^1(\mathbf{w}) + \nabla L^1(\mathbf{w})^\top(\mathbf{v} - \mathbf{w}) + \frac{\mu}{2} \|\mathbf{v} - \mathbf{w}\|^2 \quad (49)$$

for all  $\mathbf{v} \in \mathcal{W} \cap \mathcal{V}$  where  $\mathcal{V} = \{\mathbf{v} : \mathbf{v} - \mathbf{w} \in \text{span}(\mathbf{X})\}$ . Taking  $\min_{\mathbf{v} \in \mathcal{W} \cap \mathcal{V}}$  on both sides, then relaxing part of the constraint on the right-hand side yields

$$\begin{aligned} &\min_{\mathbf{v} \in \mathcal{W} \cap \mathcal{V}} L^1(\mathbf{v}) \\ &\geq \min_{\mathbf{v} \in \mathcal{W} \cap \mathcal{V}} L^1(\mathbf{w}) + \nabla L^1(\mathbf{w})^\top(\mathbf{v} - \mathbf{w}) + \frac{\mu}{2} \|\mathbf{v} - \mathbf{w}\|^2 \\ &\geq \min_{\mathbf{v} \in \mathcal{V}} L^1(\mathbf{w}) + \nabla L^1(\mathbf{w})^\top(\mathbf{v} - \mathbf{w}) + \frac{\mu}{2} \|\mathbf{v} - \mathbf{w}\|^2. \end{aligned} \quad (50)$$

Now, the minimum on the left-hand side is equal to 0 and is attained at  $\mathbf{v} = \mathbf{w} + \mathbf{P}_{\mathbf{X}}(\mathbf{w}_* - \mathbf{w})$ , as can be seen from Theorem A.2. For the right-hand side, we can substitute  $\mathbf{v} = \mathbf{w} + \mathbf{X}\mathbf{a}$  for  $\mathbf{a} \in \mathbb{R}^n$  and find the unconstrained minimum with respect to  $\mathbf{a}$ . We get

$$\begin{aligned} 0 &\geq L^1(\mathbf{w}) - \frac{1}{2\mu} \nabla L^1(\mathbf{w})^\top \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1} \mathbf{X}^\top \nabla L^1(\mathbf{w}) \\ &\geq L^1(\mathbf{w}) - \frac{\lambda_{\max}}{2\mu} \|\nabla L^1(\mathbf{w})\|^2, \end{aligned} \quad (51)$$

where  $\lambda_{\max} > 0$  is the largest eigenvalue of  $\mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top$ . This yields the result with  $c = \mu/\lambda_{\max}$ .  $\square$

## B. Proof of Theorem 1

We will prove a supporting lemma, and then the theorem.

**Lemma B.1.** *Assume the student is a directly parameterised linear classifier ( $N = 1$ ) initialised at zero,  $\mathbf{w}(0) = \mathbf{0}$ . Then,  $\mathbf{w}(\tau) \in \text{span}(\mathbf{X})$  for  $\tau \in [0, \infty)$ .*

*Proof.* Let  $\mathbf{q} \in \mathbb{R}^d$  be any vector orthogonal to the span of  $\mathbf{X}$ . It suffices to show that  $\mathbf{q}^\top \mathbf{w}(\tau) = 0$ . For that, notice that  $\mathbf{q}^\top \mathbf{w}(0) = 0$  and

$$\frac{d}{d\tau}(\mathbf{q}^\top \mathbf{w}(\tau)) = -\mathbf{q}^\top \nabla L^1(\mathbf{w}(\tau)) = 0, \quad (52)$$

where the last equality follows from the fact that  $\nabla L^1(\mathbf{w}(\tau)) \in \text{span}(\mathbf{X})$  (Theorem A.1). The claim follows.  $\square$

**Theorem 1.** *Assume the student is a directly parameterised linear classifier ( $N = 1$ ) with weight vector initialised at zero,  $\mathbf{w}(0) = \mathbf{0}$ . Then, the student's weight vector fulfills almost surely*

$$\mathbf{w}(t) \rightarrow \hat{\mathbf{w}}, \quad (5)$$

for  $t \rightarrow \infty$ , with

$$\hat{\mathbf{w}} = \begin{cases} \mathbf{w}_*, & n \geq d, \\ \mathbf{X}(\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{w}_*, & n < d. \end{cases} \quad (6)$$

*Proof.* Recall the time-derivative of  $L$ ,

$$L'(\tau) = -\|\nabla L^1(\mathbf{w}(\tau))\|^2. \quad (53)$$

The data matrix  $\mathbf{X}$  is almost surely (wrt.  $\mathbf{X} \sim P_{\mathbf{X}}^n$ ) full-rank, we can therefore apply Corollary A.1 to  $\mathcal{W} = \{\mathbf{w} : L^1(\mathbf{w}) \leq L^1(\mathbf{0})\}$  and  $\mathbf{w}(\tau)$  to lower-bound the gradient norm on the right-hand side of (53). We obtain  $L'(\tau) \leq -cL(\tau)$  for some  $c > 0$  and all  $\tau \in [0, \infty)$ , or equivalently,

$$(\log L(\tau))' \leq -c. \quad (54)$$

Integrating over  $[0, t]$  yields  $L(t) \leq L(0) \cdot e^{-ct}$ , which proves global convergence in the objective:  $L(t) \rightarrow 0$  as  $t \rightarrow \infty$ .

Now invoke Theorem A.3 with  $\mathcal{W}$  as above,  $\mathbf{v} = \mathbf{w}(t)$  and  $\mathbf{w} = \hat{\mathbf{w}}$  (we know that both  $\mathbf{w}(\tau), \hat{\mathbf{w}} \in \mathcal{W} \cap \text{span}(\mathbf{X})$ , partly by Lemma B.1):

$$L(t) \geq \frac{\mu}{2} \|\mathbf{w}(t) - \hat{\mathbf{w}}\|^2. \quad (55)$$

Since  $L(t) \rightarrow 0$  as  $t \rightarrow \infty$ , the theorem follows.  $\square$

## C. Proof of Theorem 2

**Theorem 2.** *Let  $\hat{\mathbf{w}}$  be defined as in Theorem 1. Assume the student is a deep linear network, initialized such that for some  $\epsilon > 0$ ,*

$$\|\mathbf{w}(0)\| < \min \left\{ \|\hat{\mathbf{w}}\|, \epsilon^N \left( \epsilon^2 \|\hat{\mathbf{w}}\|^{-\frac{2}{N}} + \|\hat{\mathbf{w}}\|^{2-\frac{2}{N}} \right)^{-\frac{N}{2}} \right\}, \quad (11)$$

$$L^1(\mathbf{w}(0)) < L^1(\mathbf{0}), \quad (12)$$

$$\mathbf{W}_{j+1}(0)^\top \mathbf{W}_{j+1}(0) = \mathbf{W}_j(0) \mathbf{W}_j(0)^\top \quad (13)$$

for  $j = 1, \dots, N-1$ . Then, for  $n \geq d$ , student's weight vector fulfills almost surely

$$\mathbf{w}(t) \rightarrow \hat{\mathbf{w}}, \quad (14)$$

and for  $n < d$ ,

$$\|\mathbf{w}(t) - \hat{\mathbf{w}}\| \leq \epsilon, \quad (15)$$

for all  $t$  large enough.

For the proof, we will need a result by (Arora et al., 2018), which characterises the induced flow on  $\mathbf{w}(\tau)$  when running gradient descent on the component matrices  $\mathbf{W}_i$ .

**Lemma C.1** ((Arora et al., 2018, Claim 2)). *If the balancedness condition (13) holds, then*

$$\frac{\partial \mathbf{w}(\tau)}{\partial \tau} = -\|\mathbf{w}(\tau)\|^{\frac{2(N-1)}{N}} (\nabla L^1(\mathbf{w}(\tau)) + (N-1) \cdot \mathbf{P}_{\mathbf{w}(\tau)} \nabla L^1(\mathbf{w}(\tau))). \quad (56)$$

*Proof of Theorem 2.* Similarly to the case  $N=1$ , we start by looking at the time-derivative of  $L$ ,

$$\begin{aligned} L'(\tau) &= \nabla L^1(\mathbf{w}(\tau))^\top \left( \frac{\partial \mathbf{w}(\tau)}{\partial \tau} \right) \\ &= -\|\mathbf{w}(\tau)\|^{\frac{2(N-1)}{N}} \left( \|\nabla L^1(\mathbf{w}(\tau))\|^2 \right. \\ &\quad \left. + (N-1) \cdot \|\mathbf{P}_{\mathbf{w}(\tau)} \nabla L^1(\mathbf{w}(\tau))\|^2 \right) \\ &\leq -\|\mathbf{w}(\tau)\|^{\frac{2(N-1)}{N}} \cdot \|\nabla L^1(\mathbf{w}(\tau))\|^2. \end{aligned} \quad (57)$$

It is non-positive, so  $\mathbf{w}(\tau)$  stays within the  $L(0)$ -sublevel set throughout optimisation,

$$\mathbf{w}(\tau) \in \mathcal{W} = \{\mathbf{w} : L^1(\mathbf{w}) \leq L(0)\}. \quad (58)$$

Also,  $\mathcal{W}$  is convex and by Assumption (12) it does not contain  $\mathbf{0}$ . We can therefore take  $\delta > 0$  to be the distance between  $\mathcal{W}$  and  $\mathbf{0}$ , and it follows that  $\|\mathbf{w}(\tau)\| \geq \delta$  for  $\tau \in [0, \infty)$ .

Now, noting that  $\mathbf{X}$  is almost surely full-rank, apply Corollary A.1 to  $\mathcal{W}$  and  $\mathbf{w}(\tau)$  to upper-bound the right-hand side of (57),

$$L'(\tau) \leq -c\delta^{\frac{2(N-1)}{N}} L(\tau). \quad (59)$$

Letting  $\tilde{c} = c\delta^{\frac{2(N-1)}{N}}$ , we get  $(\log L(\tau))' \leq -\tilde{c}$  and consequently  $L(t) \leq L(0) \cdot e^{-\tilde{c}t}$ . This proves convergence in the objective,  $L(t) \rightarrow 0$  as  $t \rightarrow \infty$ .

To prove convergence in parameters, we decompose the 'error'  $\mathbf{w}(\tau) - \hat{\mathbf{w}}$  into orthogonal components and bound each of them separately,

$$\begin{aligned} \|\mathbf{w}(\tau) - \hat{\mathbf{w}}\|^2 &= \|\mathbf{P}_{\mathbf{X}}(\mathbf{w}(\tau) - \hat{\mathbf{w}})\|^2 \\ &\quad + \|\mathbf{P}_{\mathbf{Q}}(\mathbf{w}(\tau) - \hat{\mathbf{w}})\|^2, \end{aligned} \quad (60)$$

where the columns of  $\mathbf{Q} \in \mathbb{R}^{d \times (d-n)}$  orthogonally complement those of  $\mathbf{X}$ . If  $n \geq d$ , we simply bound the first term and disregard the second one.

To bound the first term, invoke Theorem A.3 with  $\mathcal{W}$ ,  $\mathbf{v} = \mathbf{P}_{\mathbf{X}}\mathbf{w}(\tau)$  and  $\mathbf{w} = \mathbf{P}_{\mathbf{X}}\hat{\mathbf{w}}$ . One can check that  $L^1(\mathbf{P}_{\mathbf{X}}\mathbf{u}) = L^1(\mathbf{u})$  for all  $\mathbf{u} \in \mathbb{R}^d$ , so  $\mathbf{P}_{\mathbf{X}}\mathbf{w}(\tau) \in \mathcal{W}$  and our use of the theorem is legal. We obtain

$$L(\tau) \geq \frac{\mu}{2} \|\mathbf{P}_{\mathbf{X}}(\mathbf{w}(\tau) - \hat{\mathbf{w}})\|^2. \quad (61)$$

Since  $L(\tau) \rightarrow 0$ , it follows that

$$\|\mathbf{P}_{\mathbf{X}}(\mathbf{w}(\tau) - \hat{\mathbf{w}})\|^2 \rightarrow 0 \quad (62)$$

as  $\tau \rightarrow \infty$ .

For the second term, notice that  $\hat{\mathbf{w}} \in \text{span}(\mathbf{X})$ , so  $\mathbf{P}_{\mathbf{Q}}\hat{\mathbf{w}}$  vanishes and we are left with  $\|\mathbf{P}_{\mathbf{Q}}\mathbf{w}(\tau)\|^2$ . Denote this quantity  $q(\tau)$ . Its time derivative is

$$\begin{aligned} q'(\tau) &= 2(\mathbf{P}_{\mathbf{Q}}\mathbf{w}(\tau))^\top \left( \frac{\partial \mathbf{w}(\tau)}{\partial \tau} \right) \\ &= -2\|\mathbf{w}(\tau)\|^{\frac{2(N-1)}{N}} \left( \mathbf{w}(\tau)^\top \mathbf{P}_{\mathbf{Q}} \nabla L^1(\mathbf{w}(\tau)) + \right. \\ &\quad \left. \frac{(N-1)}{\|\mathbf{w}(\tau)\|^2} \cdot \mathbf{w}(\tau)^\top \mathbf{P}_{\mathbf{Q}} \mathbf{w}(\tau) \cdot \mathbf{w}(\tau)^\top \nabla L^1(\mathbf{w}(\tau)) \right) \\ &= -2q(\tau)(N-1)\|\mathbf{w}(\tau)\|^{-2/N} \mathbf{w}(\tau)^\top \nabla L^1(\mathbf{w}(\tau)), \end{aligned} \quad (63)$$

where we have used the fact that  $\nabla L^1(\mathbf{w}(\tau)) \in \text{span}(\mathbf{X})$  (Theorem A.1) and  $\mathbf{Q}$  is orthogonal to  $\mathbf{X}$ . Rearranging, we obtain

$$\frac{d}{d\tau} \left( \frac{\log q(\tau)}{2(N-1)} \right) = -\|\mathbf{w}(\tau)\|^{-2/N} \cdot \mathbf{w}(\tau)^\top \nabla L^1(\mathbf{w}(\tau)). \quad (64)$$

It turns out that the right-hand side expression is integrable in yet another way, namely

$$\begin{aligned} \frac{d}{d\tau} \left( \frac{1}{2N} \log \|\mathbf{w}(\tau)\|^2 \right) &= \\ &= -\|\mathbf{w}(\tau)\|^{-2/N} \cdot \mathbf{w}(\tau)^\top \nabla L^1(\mathbf{w}(\tau)). \end{aligned} \quad (65)$$

Equating the two and integrating over  $[0, t]$  yields

$$\log \frac{q(t)}{q(0)} = \frac{N-1}{N} \cdot \log \frac{\|\mathbf{w}(t)\|^2}{\|\mathbf{w}(0)\|^2}, \quad (66)$$

which implies

$$\frac{q(t)}{\|\mathbf{w}(t)\|^2} \leq \left( \frac{\|\mathbf{w}(0)\|}{\|\mathbf{w}(t)\|} \right)^{2/N}, \quad (67)$$

because  $q(0) \leq \|\mathbf{w}(0)\|^2$ .

We now bound the norm of  $\mathbf{w}(t)$ . Starting from an orthogonal decomposition similar to (60) and applying (62) with (67), we get

$$\|\mathbf{w}(t)\|^2 = \|\mathbf{P}_X \mathbf{w}(t)\|^2 + \|\mathbf{P}_Q \mathbf{w}(t)\|^2$$

$$\limsup_{t \rightarrow \infty} \|\mathbf{w}(t)\|^2 \leq \|\hat{\mathbf{w}}\|^2 + \|\mathbf{w}(0)\|^{\frac{2}{N}} \limsup_{t \rightarrow \infty} \|\mathbf{w}(t)\|^{2-\frac{2}{N}}. \quad (68)$$

Denote  $\nu := \limsup_{t \rightarrow \infty} \|\mathbf{w}(t)\|$ . By the same orthogonal decomposition, we also know that  $\nu^2 \geq \limsup_{t \rightarrow \infty} \|\mathbf{P}_X \mathbf{w}(t)\|^2 = \|\hat{\mathbf{w}}\|^2 > 0$ , so we can divide both sides above by  $\nu^2$ ,

$$1 \leq \frac{\|\hat{\mathbf{w}}\|^2}{\nu^2} + \frac{\|\mathbf{w}(0)\|^{2/N}}{\nu^{2/N}} =: f(\nu). \quad (69)$$

On the right-hand side, we now have a decreasing function of  $\nu$  that goes to zero as  $\nu \rightarrow \infty$ . However, evaluated at our specific  $\nu$ , it is lower-bounded by 1, implying an implicit upper bound for  $\nu$ .

How do we find this bound? Suppose we find some constant  $K$  such that  $f(K) \leq 1$ . Then, because  $f$  is decreasing, it must be the case that  $\nu \leq K$ . One such candidate for  $K$  is

$$K = \|\hat{\mathbf{w}}\| \cdot \left( 1 - \frac{\|\mathbf{w}(0)\|^{2/N}}{\|\hat{\mathbf{w}}\|^{2/N}} \right)^{\frac{-N}{2(N-1)}}. \quad (70)$$

(Here we have used condition (11):  $\|\mathbf{w}(0)\| < \|\hat{\mathbf{w}}\|$ .) To check that indeed  $f(K) \leq 1$ , start from the inequality

$$\begin{aligned} (\|\hat{\mathbf{w}}\|/K)^{\frac{2(N-1)}{N}} + \frac{\|\mathbf{w}(0)\|^{2/N}}{\|\hat{\mathbf{w}}\|^{2/N}} &= 1 \\ &\leq \left( 1 - \frac{\|\mathbf{w}(0)\|^{2/N}}{\|\hat{\mathbf{w}}\|^{2/N}} \right)^{\frac{-1}{N-1}} = (\|\hat{\mathbf{w}}\|/K)^{-\frac{2}{N}}. \end{aligned} \quad (71)$$

Taking the leftmost and rightmost expression and multiplying by  $(\|\hat{\mathbf{w}}\|/K)^{2/N}$  yields

$$f(K) = \frac{\|\hat{\mathbf{w}}\|^2}{K^2} + \frac{\|\mathbf{w}(0)\|^{2/N}}{K^{2/N}} \leq 1. \quad (72)$$

Hence,

$$\limsup_{t \rightarrow \infty} \|\mathbf{w}(t)\| \leq \|\hat{\mathbf{w}}\| \cdot \left( 1 - \frac{\|\mathbf{w}(0)\|^{2/N}}{\|\hat{\mathbf{w}}\|^{2/N}} \right)^{\frac{-N}{2(N-1)}}. \quad (73)$$

Finally, let us turn back to our original goal of bounding  $\|\mathbf{w}(\tau) - \hat{\mathbf{w}}\|^2$ . With (60), (62), (67) and (73), we now know that

$$\limsup_{t \rightarrow \infty} \|\mathbf{w}(\tau) - \hat{\mathbf{w}}\|^2 \quad (74)$$

$$\leq \|\mathbf{w}(0)\|^{\frac{2}{N}} \|\hat{\mathbf{w}}\|^{\frac{2(N-1)}{N}} \left( 1 - \frac{\|\mathbf{w}(0)\|^{2/N}}{\|\hat{\mathbf{w}}\|^{2/N}} \right)^{-1} \quad (75)$$

$$= \frac{\|\hat{\mathbf{w}}\|^{2+2/N}}{\|\hat{\mathbf{w}}\|^{2/N} - \|\mathbf{w}(0)\|^{2/N}} - \|\hat{\mathbf{w}}\|^2. \quad (76)$$

Hence, if we initialise close enough to zero, as specified by condition (11), we can ensure that

$$\limsup_{t \rightarrow \infty} \|\mathbf{w}(\tau) - \hat{\mathbf{w}}\|^2 < \epsilon^2. \quad (77)$$

This concludes the proof.  $\square$

## D. Theorem 3 for Approximate Distillation

We extend Theorem 3 to the setting where the student learns the solution  $\hat{\mathbf{w}} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{w}_*$  only  $\epsilon$ -approximately, as is the case for deep linear networks initialised as in Theorem 2. When  $n \geq d$ , the teacher's weight vector is recovered exactly and the transfer risk is zero, even when the student is deep. The following theorem therefore only covers the case  $n < d$ .

**Theorem D.1** (Risk bound for approximate distillation).

Let  $n < d$ . For any training set  $\mathbf{X} \in \mathbb{R}^{d \times n}$ , let  $\hat{h}_{\mathbf{X}}(\mathbf{x}) = \mathbb{1}\{\hat{\mathbf{w}}_\epsilon^\top \mathbf{x} \geq 0\}$  be a linear classifier whose weight vector is  $\epsilon$ -close to the distillation solution  $\hat{\mathbf{w}}$ , i.e.  $\|\hat{\mathbf{w}}_\epsilon - \hat{\mathbf{w}}\| \leq \epsilon$ , where  $\epsilon$  is a positive constant such that  $\epsilon \leq \frac{1}{2} \|\hat{\mathbf{w}}\|$ . Define  $\delta := \sqrt{\frac{2\pi\epsilon}{\|\hat{\mathbf{w}}\|}}$ . Then, it holds for any  $\beta \in [0, \pi/2 - \delta]$  that

$$\mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}^{\otimes n}} \left[ R(\hat{h}_{\mathbf{X}} | P_{\mathbf{x}}, \mathbf{w}_*) \right] \leq p(\beta) + p(\pi/2 - \delta - \beta)^n. \quad (78)$$

The result is very similar to Theorem 3 in the main text, the only difference is the constant  $\delta$  which compensates for the imprecision in learning  $\hat{\mathbf{w}}$  by pushing the bound up (recall that  $p$  is decreasing). However, as  $\epsilon$  goes to zero, so does  $\delta$  and we recover the original bound.

For the proof, we start with a tool for controlling the angle between  $\hat{\mathbf{w}}$  and  $\hat{\mathbf{w}}_\epsilon$ . Recall that the angle is defined as

$$\alpha(\mathbf{w}, \mathbf{v}) = \cos^{-1} \left( \frac{\mathbf{w}^\top \mathbf{v}}{\|\mathbf{w}\| \cdot \|\mathbf{v}\|} \right) \quad (79)$$

for  $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ .

**Lemma D.1.** *Let  $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$  be such that  $\|\mathbf{w} - \mathbf{v}\| \leq \epsilon$ , where  $\epsilon \leq \frac{1}{2}\|\mathbf{w}\|$ . Then  $\alpha(\mathbf{w}, \mathbf{v}) \leq \sqrt{\frac{2\pi\epsilon}{\|\mathbf{w}\|}}$ .*

*Proof of Lemma D.1.* The first step is to lower-bound the inner product  $\mathbf{w}^\top \mathbf{v}$ . To that end, we expand and rearrange  $\|\mathbf{w} - \mathbf{v}\|^2 \leq \epsilon^2$  to obtain

$$2\mathbf{w}^\top \mathbf{v} \geq \|\mathbf{w}\|^2 + \|\mathbf{v}\|^2 - \epsilon^2. \quad (80)$$

Now use the triangle relation  $\|\mathbf{v}\| \geq \|\mathbf{w}\| - \epsilon$  squared to lower-bound the right-hand side of (80) and get

$$2\mathbf{w}^\top \mathbf{v} \geq 2\|\mathbf{w}\|^2 - 2\epsilon\|\mathbf{w}\|, \quad (81)$$

which implies

$$\frac{\mathbf{w}^\top \mathbf{v}}{\|\mathbf{w}\| \cdot \|\mathbf{v}\|} \geq \frac{\|\mathbf{w}\| - \epsilon}{\|\mathbf{v}\|} \geq \frac{\|\mathbf{w}\| - \epsilon}{\|\mathbf{w}\| + \epsilon} \geq 1 - \frac{2\epsilon}{\|\mathbf{w}\|}. \quad (82)$$

Thus,

$$1 - \frac{2\epsilon}{\|\mathbf{w}\|} \leq \frac{\mathbf{w}^\top \mathbf{v}}{\|\mathbf{w}\| \cdot \|\mathbf{v}\|} = \cos(\alpha(\mathbf{w}, \mathbf{v})). \quad (83)$$

The left-hand side is by assumption non-negative, so we have  $\alpha(\mathbf{w}, \mathbf{v}) \in [-\pi/2, \pi/2]$ . On this domain,

$$\cos x \leq 1 - \frac{x^2}{\pi}, \quad (84)$$

which lets us deduce

$$1 - \frac{2\epsilon}{\|\mathbf{w}\|} \leq 1 - \frac{\alpha(\mathbf{w}, \mathbf{v})^2}{\pi}. \quad (85)$$

Rearranging yields the result.  $\square$

*Proof of Theorem D.1.* We decompose the expected risk as follows:

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim P_{\mathbf{x}}^n} \left[ R(\hat{h}_{\mathbf{X}} | P_{\mathbf{x}}, \mathbf{w}_*) \right] &= \mathbb{P}_{\substack{\mathbf{x} \sim P_{\mathbf{x}}^n \\ \mathbf{x} \sim P_{\mathbf{x}}}} [\mathbf{w}_*^\top \mathbf{x} \cdot \hat{\mathbf{w}}_\epsilon^\top \mathbf{x} < 0] = \\ &= \int_{\mathbf{x}: \bar{\alpha}(\mathbf{w}_*, \mathbf{x}) \geq \beta} \mathbb{P}_{\mathbf{x} \sim P_{\mathbf{x}}^n} [\mathbf{w}_*^\top \mathbf{x} \cdot \hat{\mathbf{w}}_\epsilon^\top \mathbf{x} < 0 | \mathbf{x}] dP_{\mathbf{x}} \\ &+ \int_{\mathbf{x}: \bar{\alpha}(\mathbf{w}_*, \mathbf{x}) < \beta, \mathbf{w}_*^\top \mathbf{x} > 0} \mathbb{P}_{\mathbf{x} \sim P_{\mathbf{x}}^n} [\hat{\mathbf{w}}_\epsilon^\top \mathbf{x} < 0 | \mathbf{x}] dP_{\mathbf{x}} \\ &+ \int_{\mathbf{x}: \bar{\alpha}(\mathbf{w}_*, \mathbf{x}) < \beta, \mathbf{w}_*^\top \mathbf{x} < 0} \mathbb{P}_{\mathbf{x} \sim P_{\mathbf{x}}^n} [\hat{\mathbf{w}}_\epsilon^\top \mathbf{x} > 0 | \mathbf{x}] dP_{\mathbf{x}}. \end{aligned} \quad (86)$$

Let us fix some  $\mathbf{x}$  for which  $\bar{\alpha}(\mathbf{w}_*, \mathbf{x}) < \beta$  and  $\mathbf{w}_*^\top \mathbf{x} > 0$ ; for this  $\mathbf{x}$  we have  $\alpha(\mathbf{w}_*, \mathbf{x}) = \bar{\alpha}(\mathbf{w}_*, \mathbf{x})$ . Consider the

situation where  $\bar{\alpha}(\mathbf{w}_*, \mathbf{x}_i) < \pi/2 - \beta - \delta$  for some  $i$ . Then by the triangle inequality, Lemma D.1 and Lemma 1,

$$\alpha(\hat{\mathbf{w}}_\epsilon, \mathbf{x}) \leq \alpha(\hat{\mathbf{w}}_\epsilon, \hat{\mathbf{w}}) + \alpha(\mathbf{w}_*, \hat{\mathbf{w}}) + \alpha(\mathbf{w}_*, \mathbf{x}) \quad (87)$$

$$\leq \delta + \bar{\alpha}(\mathbf{w}_*, \mathbf{x}_i) + \bar{\alpha}(\mathbf{w}_*, \mathbf{x}) \quad (88)$$

$$< \pi/2, \quad (89)$$

which implies  $\hat{\mathbf{w}}_\epsilon^\top \mathbf{x} > 0$ , i.e. a correct prediction (same as the teacher's). Conversely, an error can occur only if  $\bar{\alpha}(\mathbf{w}_*, \mathbf{x}_i) \geq \pi/2 - \delta - \beta$  for all  $i$ . Because  $\mathbf{x}_i$  are independent, we have

$$\begin{aligned} \mathbb{P}_{\mathbf{x} \sim P_{\mathbf{x}}^n} [\hat{\mathbf{w}}_\epsilon^\top \mathbf{x} < 0 | \mathbf{x} : \bar{\alpha}(\mathbf{w}_*, \mathbf{x}) < \beta, \mathbf{w}_*^\top \mathbf{x} > 0] \\ \leq \mathbb{P}_{\mathbf{x} \sim P_{\mathbf{x}}^n} [\forall_i : \bar{\alpha}(\mathbf{w}_*, \mathbf{x}_i) \geq \pi/2 - \delta - \beta] \\ = p(\pi/2 - \delta - \beta)^n. \end{aligned} \quad (90)$$

By a symmetric argument, one can show that

$$\begin{aligned} \mathbb{P}_{\mathbf{x} \sim P_{\mathbf{x}}^n} [\hat{\mathbf{w}}_\epsilon^\top \mathbf{x} > 0 | \mathbf{x} : \bar{\alpha}(\mathbf{w}_*, \mathbf{x}) < \beta, \mathbf{w}_*^\top \mathbf{x} < 0] \\ \leq p(\pi/2 - \delta - \beta)^n. \end{aligned} \quad (91)$$

Combining (86), (90) and (91) yields the result.  $\square$