

---

# An Optimization and Generalization Analysis for Max-Pooling Networks

## Supplementary Material

---

Alon Brutzkus<sup>1</sup>

Amir Globerson<sup>1</sup>

<sup>1</sup>The Blavatnik School of Computer Science, Tel Aviv University

### A CONVERGENCE RATES FOR THEOREM 5.1

In Ji and Telgarsky [2019], Theorem 4.2, they show the following for logistic regression initialized at zero and a certain learning rate schedule. The margin of the learned classifier is  $\frac{\gamma}{2}$  where  $\gamma$  is the max-margin after  $O\left(\frac{1}{\gamma^2}\right)$  iterations.<sup>1</sup> They show this for normalized points with norm 1. In our case (see the proof of Theorem 5.1), the max margin after normalizing the points to have norm 1, is  $\frac{1}{\sqrt{d}}$ . Thus, under their assumptions, after  $O(d)$  iterations we converge to a solution whose margin is a  $\frac{1}{2}$ -multiplicative approximation of the max margin. Therefore, we obtain for this solution, up to a constant, the same generalization guarantees as the max margin classifier (which we provide in the theorem).

### B PROOF OF LEMMA 5.2

By definition of the initialization we have  $\mathbb{P}(i \in \mathcal{A}^+) = \frac{1}{2}$ . Furthermore, we have that  $\mathbb{P}(i \in \mathcal{W}_0^+) = \frac{(1-2^{-d+1})}{d-1}$ . This follows, since with probability  $2^{-d+1}$ , for all  $\mathbf{o} \in \mathcal{O} \setminus \{2\}$ ,  $\mathbf{w}_i^{(0)} \cdot \mathbf{o} \leq 0$ . On the other hand, with probability  $(1 - 2^{-d+1})$ , there exists at least one  $\mathbf{o} \in \mathcal{O} \setminus \{2\}$  such that  $\mathbf{w}_i^{(0)} \cdot \mathbf{o} > 0$ . Assume we condition on the latter event. Then, we get by symmetry that  $\mathbf{o}_1$  maximizes the dot product with  $\mathbf{w}_i^{(0)}$ , among patterns in  $\mathcal{O} \setminus \{2\}$ , with probability  $\frac{1}{d-1}$ .

By independence of  $W_0$  and  $\mathbf{a}^{(0)}$ , we have:  $\mathbb{P}(i \in \mathcal{W}_0^+ \cap \mathcal{A}^+) = \frac{(1-2^{-d+1})}{2(d-1)}$ . Then, by Hoeffding's inequality we get:

$$\mathbb{P}\left(\left|\frac{|\mathcal{W}_0^+ \cap \mathcal{A}^+|}{k} - \frac{(1-2^{-d+1})}{2(d-1)}\right| > \frac{1}{4d}\right) \leq 2e^{-2k(\frac{1}{4d})^2} \leq 2e^{-d} \quad (1)$$

where in the last inequality we used the assumption on  $k$ . Since  $\frac{(1-2^{-d+1})}{2(d-1)} \geq \frac{1}{2d}$  and  $\frac{(1-2^{-d+1})}{2(d-1)} \leq \frac{1}{d}$  for  $d \geq 3$ , we get that with probability at least  $1 - 2e^{-d}$ ,  $|\mathcal{W}_0^+ \cap \mathcal{A}^+| \geq \frac{(1-2^{-d+1})k}{2(d-1)} - \frac{k}{4d} \geq \frac{k}{4d}$  and  $|\mathcal{W}_0^+ \cap \mathcal{A}^+| \leq \frac{(1-2^{-d+1})k}{2(d-1)} + \frac{k}{4d} \leq \frac{k}{d}$ . By the symmetry of our problem and definitions of the sets  $\mathcal{W}_0^+$ ,  $\mathcal{W}_0^-$ ,  $\mathcal{A}^+$ ,  $\mathcal{A}^-$ , we similarly get that with probability at least  $1 - 2e^{-d}$ ,  $\frac{k}{4d} \leq |\mathcal{W}_0^- \cap \mathcal{A}^-| \leq \frac{k}{d}$ . Applying the union bound concludes the proof.

### C PROOF OF LEMMA 5.3

We first prove the following two auxiliary lemmas.

**Lemma C.1.** *For all  $0 \leq t \leq T_1$  and all  $1 \leq i \leq k$ ,  $\|\mathbf{w}_i^{(t)}\| \leq \eta_1(t+1)$ .*

<sup>1</sup> $\mathcal{O}$  hides a dependency on  $\log m$ .

*Proof.* First we notice that for all  $1 \leq i \leq k$ ,  $\left\| \frac{\partial \mathcal{L}_1}{\partial \mathbf{w}_i} (W, \mathbf{a}^{(0)}) \right\| \leq 1$ . This follows since for all  $1 \leq j \leq n$  and all  $\mathbf{x} \in S_1$ ,  $\|\mathbf{x}[j]\| = 1$  (recall that  $\|\mathbf{o}\| = 1$  for  $\mathbf{o} \in \mathcal{O}$ ).

Therefore, for all  $0 \leq t \leq T_1$  and  $1 \leq i \leq k$ ,  $\|\mathbf{w}_i^{(t)}\| \leq r + \eta_1 t \leq \eta_1(t+1)$ .  $\square$

**Lemma C.2.** For all  $\mathbf{x} \in S_1$  and  $0 \leq t \leq T_1$   $|N_{\text{CNN}}(\mathbf{x}; (W^{(t)}, a^{(0)}))| \leq \frac{1}{2}$ .

*Proof.* By Lemma C.1 we have for all  $\mathbf{x} \in S_1$ :

$$\begin{aligned} |N_{\text{CNN}}(\mathbf{x}; (W^{(t)}, a^{(0)}))| &= \left| \sum_{i=1}^k a_i^{(0)} \left[ \max_j \left\{ \sigma(\mathbf{w}_i^{(t)} \cdot \mathbf{x}[j]) \right\} \right] \right| \\ &\leq k \max_{1 \leq i \leq k} \|\mathbf{w}_i^{(t)}\| \max_{1 \leq j \leq n} \|\mathbf{x}[j]\| \\ &\leq k \eta_1 (t+1) \\ &\leq \frac{1}{2} \end{aligned}$$

where the last inequality follows by the assumption on  $\eta_1$ .  $\square$

Lemma 5.3 follows by the following lemma.

**Lemma C.3.** With probability at least  $1 - 4e^{-\frac{m}{36}}$ , for all  $0 \leq t \leq T_1$  and all  $i \in \mathcal{W}_0^+ \cap \mathcal{A}^+$  the following holds:

1.  $\mathbf{o}_1 \cdot \mathbf{w}_i^{(t)} \geq \frac{t\eta_1}{9}$ .
2. For all  $j \neq 1$ , it holds that  $\mathbf{o}_j \cdot \mathbf{w}_i^{(t)} \leq r$ .

*Proof.* We will prove the claim for  $i \in \mathcal{W}_0^+ \cap \mathcal{A}^+$ . We prove the two claims by induction on  $t$ . In the proof by induction we also show a third claim that: for all  $\mathbf{x}_+ \in S_1^+$ ,  $\mathbf{p}_t^{(i)}(\mathbf{x}_+) = \mathbf{o}_1$ .

For the proof, we condition on the event:

$$\frac{|S_1^+|}{m_1}, \frac{|S_1^-|}{m_1} \geq \frac{m_1}{3} \quad (2)$$

This holds with probability at least  $1 - 4e^{-\frac{m}{36}}$  by applying Hoeffding's inequality and a union bound (over positive and negative samples).

For  $t = 0$ , we have by definition for all  $i \in \mathcal{W}_0^+ \cap \mathcal{A}^+$ ,  $\mathbf{o}_1 \cdot \mathbf{w}_i^{(0)} > 0$ . The second claim holds by the definition of the initialization. The third claim follows by the definition of  $\mathcal{W}_0^+ \cap \mathcal{A}^+$ .

Assume the three claims above hold for  $t = T$ . We will prove them for  $t = T + 1$ .

Proof of Claim 1. By the gradient update in the first layer, the following holds for  $i \in \mathcal{W}_0^+ \cap \mathcal{A}^+$ :

$$\begin{aligned} \mathbf{w}_i^{(T+1)} &= \mathbf{w}_i^{(T)} - \frac{\eta_1}{m_1} \sum_{\mathbf{x}_+ \in S_1^+} \ell' (N_{\text{CNN}}(\mathbf{x}_+; (W^{(T)}, a^{(0)}))) \mathbf{p}_T^{(i)}(\mathbf{x}_+) \\ &\quad + \frac{\eta_1}{m_1} \sum_{\mathbf{x}_- \in S_1^-} \ell' (-N_{\text{CNN}}(\mathbf{x}_-; (W^{(T)}, a^{(0)}))) \mathbf{p}_T^{(i)}(\mathbf{x}_-) \end{aligned} \quad (3)$$

where  $\ell'(z) = -\frac{1}{1+e^z}$  is the derivative of the logistic loss. Note that for all  $z$ ,  $|\ell'(z)| \leq 1$ . Therefore, for all  $\mathbf{x}_- \in S_1^-$ , we have:

$$|\ell'(-N_{\text{CNN}}(\mathbf{x}_-; (W^{(T)}, a^{(0)})))| \leq 1 \quad (4)$$

By Lemma C.2 we have for all  $\mathbf{x} \in S_1$   $|N_{\text{CNN}}(\mathbf{x}; (W^{(T)}, a^{(0)}))| \leq \frac{1}{2}$ . Therefore, for all  $\mathbf{x}_+ \in S_1^+$ :

$$|\ell'(N_{\text{CNN}}(\mathbf{x}_+; (W^{(T)}, a^{(0)})))| \geq \frac{1}{1 + \sqrt{e}} \geq \frac{1}{3} \quad (5)$$

By the induction hypothesis, we have for  $i \in \mathcal{W}_0^+ \cap \mathcal{A}^+$  and all  $\mathbf{x}_+ \in S_1^+$  that  $\mathbf{p}_T^{(i)}(\mathbf{x}_+) = \mathbf{o}_1$ . Therefore we have:

$$\mathbf{p}_T^{(i)}(\mathbf{x}_+) \cdot \mathbf{o}_1 = 1 \quad (6)$$

For all  $\mathbf{x}_- \in S_1^-$ , we have  $\mathbf{p}_T^{(i)}(\mathbf{x}_-) = \mathbf{o}_j$  for  $j \neq 1$  that depends on  $\mathbf{x}_-$ . Therefore:

$$\mathbf{p}_T^{(i)}(\mathbf{x}_-) \cdot \mathbf{o}_1 = 0 \quad (7)$$

By the facts above we complete the proof of the first claim:

$$\begin{aligned} \mathbf{w}_i^{(T+1)} \cdot \mathbf{o}_1 &\stackrel{\text{Eq. 3,4,5}}{\geq} \mathbf{w}_i^{(T)} \cdot \mathbf{o}_1 + \frac{\eta_1}{3m_1} \sum_{\mathbf{x}_+ \in S_1^+} \mathbf{p}_T^{(i)}(\mathbf{x}_+) \cdot \mathbf{o}_1 \\ &\quad - \frac{\eta_1}{m_1} \sum_{\mathbf{x}_- \in S_1^-} \mathbf{p}_T^{(i)}(\mathbf{x}_-) \cdot \mathbf{o}_1 \\ &\stackrel{\text{Eq. 2,6,7}}{\geq} \mathbf{w}_i^{(T)} \cdot \mathbf{o}_1 + \frac{\eta_1}{9} \\ &\geq \frac{(T+1)\eta_1}{9} \end{aligned} \quad (8)$$

where the last inequality follows from the induction hypothesis.

Proof of Claim 2. Since for all  $\mathbf{x}_+ \in S_1^+$ ,  $\mathbf{p}_T^{(i)}(\mathbf{x}_+) = \mathbf{o}_1$  we have for all  $1 \leq j \leq d$ ,  $j \neq 1$ :

$$\mathbf{p}_T^{(i)}(\mathbf{x}_+) \cdot \mathbf{o}_j = 0 \quad (9)$$

By the facts (1) for all  $\mathbf{x}_- \in S_1^-$  and  $j \neq 1$  it holds that  $\mathbf{p}_T^{(i)}(\mathbf{x}_-) \cdot \mathbf{o}_j \geq 0$  and (2)  $l'(z) < 0$  for all  $z$ , we have:

$$\frac{\eta_1}{m_1} \sum_{\mathbf{x}_- \in S_1^-} \ell'(-N_{\text{CNN}}((; (W^{(T)}, a^{(0)})) \mathbf{x}_-)) \mathbf{p}_T^{(i)}(\mathbf{x}_-) \cdot \mathbf{o}_j \leq 0 \quad (10)$$

Therefore we have for  $j \neq 1$ :

$$\mathbf{w}_i^{(T+1)} \cdot \mathbf{o}_j \stackrel{\text{Eq.9,10}}{\leq} \mathbf{w}_i^{(T)} \cdot \mathbf{o}_j \leq r \quad (11)$$

where the right inequality follows by the induction hypothesis.

Proof of Claim 3. Since  $r < \frac{\eta_1(T+1)}{9}$  we conclude by Eq. 8 and Eq. 11 that for all  $\mathbf{x}_+ \in S_1^+$ ,  $\mathbf{p}_{T+1}^{(i)}(\mathbf{x}_+) = \mathbf{o}_1$ .  $\square$

## D PROOF OF LEMMA 5.5

By Lemma C.1, for all  $1 \leq t \leq T_1$  and  $1 \leq i \leq k$ ,  $\|\mathbf{w}_i^{(t)}\| \leq \eta_1(t+1)$ . Therefore, for all  $1 \leq j \leq d$  and  $\mathbf{x}$  sampled from  $\mathcal{D}$ ,  $\mathbf{x}[j] \cdot \mathbf{w}_i^{(t)} \leq 2\eta_1 t$ .

## E PROOF OF PART 3 OF THEOREM 5.1

Here we condition on the events of previous lemmas which hold with probability at least  $1 - 4e^{-d} - 4e^{-\frac{m}{36}}$ . For each  $\mathbf{x}$  sampled from  $\mathcal{D}$ , define  $\mathbf{z}(\mathbf{x}) \in \mathbb{R}^k$  such that for all  $1 \leq i \leq k$ , its  $i$ th entry is  $z_i(\mathbf{x}) = \max_j \left\{ \sigma \left( \mathbf{w}_i^{(T_1)} \cdot \mathbf{x}[j] \right) \right\}$ . Notice that by Eq. 3 we have  $z_i(\mathbf{x}) = \mathbf{w}_i^{(T_1)} \cdot \mathbf{p}_i^{(T_1)}(\mathbf{x})$ . Define a new distribution of points  $\mathcal{D}_z$  over  $\mathbb{R}^k \times \{\pm 1\}$ , which samples a point  $(\mathbf{z}(\mathbf{x}), y)$  where  $(\mathbf{x}, y) \sim \mathcal{D}$ .

Our goal is to show that  $\mathcal{D}_z$  is linearly separable and can be separated with a classifier of relatively low norm. Then, we will use recent results on logistic regression, which show that GD converges to low norm solutions. Therefore, by optimizing the

second layer,  $\text{LW}_{\text{CNN}}$  will converge to a low norm solution. Finally, we will apply norm-based generalization bounds to obtain a generalization guarantee for  $\text{LW}_{\text{CNN}}$ .

First we will show that  $\mathcal{D}_z$  is linearly separable. Indeed define  $\mathbf{v}^* \in \mathbb{R}^k$  as follows. For  $i \in \mathcal{W}_0^+ \cap \mathcal{A}^+$  let  $\mathbf{v}_i^* = \frac{80d}{k\eta_1 T_1}$  and for  $i \in \mathcal{W}_0^- \cap \mathcal{A}^-$  let  $\mathbf{v}_i^* = -\frac{80d}{k\eta_1 T_1}$ . Set all other entries of  $\mathbf{v}^*$  to 0. Then for any  $\mathbf{z}(\mathbf{x}_+)$  such that  $(\mathbf{x}_+, 1) \sim \mathcal{D}$ , we have:

$$\begin{aligned} \mathbf{z}(\mathbf{x}_+) \cdot \mathbf{v}^* &= \frac{80d}{k\eta_1 T_1} \sum_{i \in \mathcal{W}_0^+ \cap \mathcal{A}^+} \mathbf{w}_i^{(T_1)} \cdot \mathbf{p}_i^{(T_1)}(\mathbf{x}_+) \\ &\quad - \frac{80d}{k\eta_1 T_1} \sum_{i \in \mathcal{W}_0^- \cap \mathcal{A}^-} \mathbf{w}_i^{(T_1)} \cdot \mathbf{p}_i^{(T_1)}(\mathbf{x}_+) \\ &> \left( \frac{80d}{k\eta_1 T_1} \right) \left( \frac{k}{4d} \right) \left( \frac{\eta_1 T_1}{10} \right) \\ &\quad - \left( \frac{80d}{k\eta_1 T_1} \right) \left( \frac{k}{d} \right) \left( \frac{\eta_1 T_1}{80} \right) \\ &= 1 \end{aligned}$$

where the inequality follows by Lemma 5.2, Lemma 5.3 and Corollary 5.4. By symmetry, we have  $-\mathbf{z}(\mathbf{x}_-) \cdot \mathbf{v}^* > 1$  for all  $(\mathbf{x}_-, -1) \sim \mathcal{D}$ .

Next, we proceed to apply Theorem 3 in Soudry et al. [2018]. It requires that  $\eta_2 < 2\beta^{-1}\sigma_{\max}^{-2}(Z)m_2^2$  where  $\beta$  is the smoothness parameter of the logistic loss,  $Z \in \mathbb{R}^{k \times m_2}$  is the matrix which contains  $\mathbf{z}(\mathbf{x}_{i+\lceil \frac{m}{2} \rceil})$  in its  $i$ th column and  $\sigma_{\max}(Z)$  is the maximum singular value of  $Z$ . In our setting,  $\beta = 1$  and by Lemma 5.5  $\sigma_{\max}^2(Z) \leq \|Z\|_F^2 \leq 4m_2 k \eta_1^2 T_1^2 \leq \frac{m_2}{4k}$ . Thus, by our assumption  $\eta_2 < 8k \leq 2\sigma_{\max}^{-2}(Z)m_2^2$  holds.

Therefore, by this theorem we are guaranteed that:

$$\lim_{t \rightarrow \infty} \frac{\mathbf{a}^{(t)}}{\|\mathbf{a}^{(t)}\|} = \frac{\hat{\mathbf{a}}}{\|\hat{\mathbf{a}}\|} \quad (12)$$

where

$$\hat{\mathbf{a}} = \arg \min_{\mathbf{v} \in \mathbb{R}^k} \|\mathbf{v}\|^2 \quad \text{s.t.} \quad \forall i \quad y_i \mathbf{v} \cdot \mathbf{z}(\mathbf{x}_i) \geq 1 \quad (13)$$

Specifically, gradient descent converges to zero training loss, i.e.,  $\lim_{T_2 \rightarrow \infty} \mathcal{L}_2((W_{T_1}, \mathbf{a}_{T_2})) = 0$ .

By optimality of  $\hat{\mathbf{a}}$  and Lemma 5.2 we have  $\|\hat{\mathbf{a}}\|^2 \leq \|\mathbf{v}^*\|^2 \leq \frac{80^2 d^2}{k^2 \eta_1^2 T_1^2} \frac{2k}{d} = \frac{2 \cdot 80^2 d}{k \eta_1^2 T_1^2}$ . Furthermore,  $\|\mathbf{z}(\mathbf{x})\|^2 \leq 4k \eta_1^2 T_1^2$  by Lemma 5.5. Therefore, we have  $\|\hat{\mathbf{a}}\|^2 \|\mathbf{z}(\mathbf{x})\|^2 = O(d)$ . Thus, by a standard margin generalization bound (e.g. Theorem 26.13 in Shalev-Shwartz and Ben-David [2014] or Bartlett and Mendelson [2002]) we have with probability at least  $1 - \delta$ :

$$\begin{aligned} &\lim_{T_2 \rightarrow \infty} \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left( \text{sign} \left( N_{\text{CNN}}(\mathbf{x}; (W^{(T_1)}, \mathbf{a}^{(T_2)})) \right) \neq y \right) \\ &= \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}} \left( \text{sign} \left( N_{\text{CNN}} \left( \mathbf{x}; \left( W^{(T_1)}, \frac{\hat{\mathbf{a}}}{\|\hat{\mathbf{a}}\|} \right) \right) \right) \neq y \right) \\ &= O \left( \sqrt{\frac{d}{m}} \right) \end{aligned}$$

where  $O$  hides an additive term which depends on  $\delta$ .

## References

Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.

Ziwei Ji and Matus Telgarsky. A refined primal-dual analysis of the implicit bias. *arXiv preprint arXiv:1906.04540*, 2019.

<sup>2</sup>We added the factor  $m_2$  because Soudry et al. [2018] consider the empirical loss without dividing by the number of samples.

Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.

Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *The Journal of Machine Learning Research*, 19(1):2822–2878, 2018.