
Appendix [Margins, Kernels and Non-linear Smoothed Perceptrons]

1. Unified Proof By Induction of Lemma 5, 8: $L_{\mu_k}(\alpha_k) \leq -\frac{1}{2}\|p_k\|_G^2$

Let $d(p)$ be 1-strongly convex with respect to the $\#$ -norm, ie $d(q) - d(p) - \langle \nabla d(p), q - p \rangle \geq \frac{1}{2}\|q - p\|_\#^2$ for any $p, q \in \Delta_n$. Let the $\#$ -norm be lower bounded by the G-norm as $\|p\|_G^2 \leq \lambda_\# \|p\|_\#^2$. For $d(p) = \sum_i p_i \log p_i + \log n$, $\#$ is the 1-norm, $\lambda_\# = 1$ and $p^* = \frac{1}{n}$. For $d(p) = \frac{1}{2}\|q - p\|_2^2$, $\#$ is the 2-norm, $\lambda_\# = n$ and $p^* = q$. Choose $\mu_0 = 2\lambda_\#$.

Let the smoothed minimizer be defined by $p_\mu(\alpha) := \arg \min_{p \in \Delta_n} \langle G\alpha, p \rangle + \mu d(p)$, and $p^* := \arg \min_{p \in \Delta_n} d(p)$. The optimality condition of $p_\mu(\alpha)$ and p^* (the gradient is perpendicular to any feasible direction) is that for any $r \in \Delta_n$,

$$\langle G\alpha + \mu \nabla d(p_\mu(\alpha)), r - p \rangle = 0 \quad (1)$$

$$\langle \nabla d(p^*), r - p \rangle = 0 \Rightarrow d(p_0) \geq \frac{1}{2}\|p_0 - p^*\|_\#^2. \quad (2)$$

$$\begin{aligned} \text{For } k=0 : \quad -\frac{1}{2}\|p_0\|_G^2 &= -\frac{1}{2}\|p_0 - p^*\|_G^2 - \langle p^*, p_0 - p^* \rangle_G - \frac{1}{2}\|p^*\|_G^2 \quad \text{writing } p_0 = (p_0 - p^*) + p^* \\ &\geq -\frac{\lambda_\#}{2}\|p_0 - p^*\|_\#^2 - \langle p^*, p_0 \rangle_G + \frac{1}{2}\|p^*\|_G^2 \quad \text{using } \|p\|_G^2 \leq \lambda_\# \|p\|_\#^2 \\ &\geq -\mu_0 d(p_0) - \langle \alpha_0, p_0 \rangle_G + \frac{1}{2}\|\alpha_0\|_G^2 \quad \text{adding } -\frac{\lambda_\#}{2}\|p_0 - p^*\|_1^2, \text{ using Eq. (2)} \\ &= L_{\mu_0}(\alpha_0). \end{aligned}$$

Assume it holds upto k . We drop index k , and write x_+ for x_{k+1} . Let $\hat{p} = (1-\theta)p + \theta p_\mu(\alpha)$ so $\alpha_+ = (1-\theta)\alpha + \theta\hat{p}$. (3)

$$\begin{aligned} L_{\mu_+}(\alpha_+) &= \frac{1}{2}\|\alpha_+\|_G^2 - \left\langle \alpha_+, p_{\mu_+}(\alpha_+) \right\rangle_G - \mu_+ d(p_{\mu_+}(\alpha_+)) \\ &= \frac{1}{2}\|(1-\theta)\alpha + \theta\hat{p}\|_G^2 - \theta \left\langle \hat{p}, p_{\mu_+}(\alpha_+) \right\rangle_G - (1-\theta) \left[\left\langle \alpha, p_{\mu_+}(\alpha_+) \right\rangle_G + \mu d(p_{\mu_+}(\alpha_+)) \right] \quad \text{using Eq. (3)} \\ &\leq (1-\theta) \left[\frac{1}{2}\|\alpha\|_G^2 - \left\langle \alpha, p_{\mu_+}(\alpha_+) \right\rangle_G - \mu d(p_{\mu_+}(\alpha_+)) \right]_1 + \theta \left[-\frac{1}{2}\|\hat{p}\|_G^2 - \left\langle \hat{p}, p_{\mu_+}(\alpha_+) - \hat{p} \right\rangle_G \right], \end{aligned}$$

where we used the convexity of $\|\cdot\|_G^2$. Recall $p_+ = (1-\theta)p + \theta p_{\mu_+}(\alpha_+)$, so that $p_+ - \hat{p} = \theta(p_{\mu_+}(\alpha_+) - p_\mu(\alpha))$. (4)

$$\begin{aligned} [\cdot]_1 &= \left[\frac{1}{2}\|\alpha\|_G^2 - \left\langle \alpha, p_\mu(\alpha) \right\rangle_G - \mu d(p_\mu(\alpha)) \right] - \left\langle \alpha, p_{\mu_+}(\alpha_+) - p_\mu(\alpha) \right\rangle_G - \mu \left[d(p_{\mu_+}(\alpha_+)) - d(p_\mu(\alpha)) \right] \\ &= L_\mu(\alpha) - \mu \left[d(p_{\mu_+}(\alpha_+)) - d(p_\mu(\alpha)) - \left\langle \nabla d(p_\mu(\alpha)), p_{\mu_+}(\alpha_+) - p_\mu(\alpha) \right\rangle \right] \quad \text{using Eq. (1)} \\ &\leq -\frac{1}{2}\|p\|_G^2 - \frac{\mu}{2}\|p_{\mu_+}(\alpha_+) - p_\mu(\alpha)\|_\#^2 \quad \text{using strong convexity of } d(p) \\ &\leq -\frac{1}{2}\|\hat{p} + (p - \hat{p})\|_G^2 - \frac{\mu}{2\lambda_\#} \|p_{\mu_+}(\alpha_+) - p_\mu(\alpha)\|_G^2 \quad \text{using } \|p\|_G^2 \leq \lambda_\# \|p\|_\#^2 \\ &\leq -\frac{1}{2}\|\hat{p}\|_G^2 - \left\langle \hat{p}, p - \hat{p} \right\rangle_G - \frac{\mu}{2\lambda_\# \theta^2} \|p_+ - \hat{p}\|_G^2 \quad \text{using Eq. (4) and dropping a } -\frac{1}{2}\|p - \hat{p}\|_G^2 \text{ term.} \end{aligned}$$

Using $(1-\theta)(p - \hat{p}) = -\theta(p_\mu(\alpha) - \hat{p})$ and substituting back,

$$\begin{aligned} L_{\mu_+}(\alpha_+) &\leq (1-\theta) \left[-\frac{1}{2}\|\hat{p}\|_G^2 + \frac{\theta}{1-\theta} \left\langle \hat{p}, p_\mu(\alpha) - \hat{p} \right\rangle_G - \frac{\mu}{2\lambda_\# \theta^2} \|p_+ - \hat{p}\|_G^2 \right] + \theta \left[-\frac{1}{2}\|\hat{p}\|_G^2 - \left\langle \hat{p}, p_{\mu_+}(\alpha_+) - \hat{p} \right\rangle_G \right] \\ &= -\frac{1}{2}\|\hat{p}\|_G^2 - \theta \left\langle \hat{p}, p_{\mu_+}(\alpha_+) - p_\mu(\alpha) \right\rangle_G - \frac{\mu(1-\theta)}{2\lambda_\# \theta^2} \|p_+ - \hat{p}\|_G^2 \\ &\leq -\frac{1}{2}\|\hat{p}\|_G^2 - \left\langle \hat{p}, p_+ - \hat{p} \right\rangle_G - \frac{1}{2}\|p_+ - \hat{p}\|_G^2 \quad \text{using Eq. (4) and } \frac{\theta^2}{1-\theta} = \frac{4}{(k+1)(k+3)} \leq \frac{4}{(k+1)(k+2)} = \frac{\mu}{\lambda_\#} \\ &= -\frac{1}{2}\|p_+\|_G^2. \end{aligned}$$

This wraps up our unified proof for both settings.