

#### B.4 APPROXIMATION OF $\tilde{\mathcal{L}}(\hat{W}, t)$ IN EQ. 22.

**Lemma B.1.** *Let  $\mathcal{L}_T(\hat{W})$  be the training loss function on which we apply gradient descent. The  $t^{\text{th}}$  iterate of gradient descent matches the minimum of  $\tilde{\mathcal{L}}_T(\hat{W}, t)$  defined as,*

$$\tilde{\mathcal{L}}_T(\hat{W}, t) := \frac{1}{2n} \sum [\hat{y}^\mu - y^\mu]^2 + \frac{1}{\eta t} \|\hat{W}\|_2^2. \quad (92)$$

*Proof.* The goal is to show,

$$\hat{W}_t = \arg \min_{\hat{W}} \tilde{\mathcal{L}}_T(\hat{W}, t), \quad \text{where,} \quad \hat{W}_t := \hat{W}_{t-1} - \eta \nabla_{\hat{W}_{t-1}} \mathcal{L}(\hat{W}_{t-1}). \quad (93)$$

For brevity of derivations, here we only consider the case where  $\lambda = \sigma_\epsilon^2 = 0$ . Recall the closed-form derivation of  $\hat{W}_t$  in Eq. 18,

$$\hat{W}_t = \left( I - [I - \eta X^T X]^t \right) (X^T X)^{-1} X^T y, \quad (94)$$

$$= \arg \min_{\hat{W}} \left[ X \hat{W} - X \left( I - [I - \eta X^T X]^t \right) (X^T X)^{-1} X^T y \right]^2, \quad (95)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} \left( I - [I - \eta X^T X]^t \right) \underbrace{(X^T X)^{-1} X^T y}_{=W, \text{ assuming } \sigma_\epsilon^2=0} \right]^2, \quad (96)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - \underbrace{x^{\mu T} \left( I - [I - \eta X^T X]^t \right) W}_{\text{a dynamic target (function of } t)} \right]^2, \quad (97)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} V \left( I - [I - \eta \Lambda]^t \right) V^T W \right]^2, \quad (X^T X = V \Lambda V^T) \quad (98)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} V \left( I - \exp(t \log[I - \eta \Lambda]) \right) V^T W \right]^2, \quad (99)$$

$$\approx \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} V \left( I - \exp(-\eta \Lambda t) \right) V^T W \right]^2, \quad (\log(1+x) \approx x) \quad (100)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} V \left( I - \exp\left(-\frac{\Lambda}{1/\eta t}\right) \right) V^T W \right]^2, \quad (101)$$

$$\approx \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} V \left( I - \exp\left(-\log\left(\frac{\Lambda}{1/\eta t} + I\right)\right) \right) V^T W \right]^2, \quad (\log(1+x) \approx x) \quad (102)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} V \left( I - \left[ \Lambda + \frac{1}{\eta t} I \right]^{-1} \frac{1}{\eta t} \right) V^T W \right]^2, \quad (103)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} V \left( \left( \Lambda + \frac{1}{\eta t} I \right)^{-1} \Lambda \right) V^T W \right]^2, \quad (104)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} \left( X^T X + \frac{1}{\eta t} I \right)^{-1} X^T X W \right]^2, \quad (105)$$

$$= \arg \min_{\hat{W}} \frac{1}{2n} \sum \left[ \hat{y}^\mu - x^{\mu T} \underbrace{\left( X^T X + \frac{1}{\eta t} I \right)^{-1} X^T y}_{\text{the normal equation}} \right]^2, \quad (106)$$

$$= \arg \min_{\hat{W}} \underbrace{\frac{1}{2n} \sum [\hat{y}^\mu - y^\mu]^2}_{\tilde{\mathcal{L}}_T(\hat{W}, t)} + \frac{1}{\eta t} \|\hat{W}\|_2^2, \quad (107)$$

which concludes the proof. Consistent with (Ali et al., 2019), this approximation implies that "L2 regularization and early stopping can be seen as equivalent (at least under the quadratic approximation of the objective function)." (Goodfellow et al., 2016).