

Self-teaching for machine translation

Tianyi Zhou

November 14, 2021

$$\begin{aligned} \min_A \quad & L(D_{val}; V^*(W^*(A))) \\ \text{s.t} \quad & V^*(W^*(A)) = \operatorname{argmin}_V \sum_{i=1}^M l(u_i, g(u_i; W^*(A)); V) \\ & W^*(A) = \operatorname{argmin}_W \sum_{i=1}^N a_i l(x_i, y_i; W) \end{aligned}$$

1 Optimization Algorithm

One-step gradient decent of $\widetilde{W}(A)$:

$$\widetilde{W}(A) \doteq \overline{W}(A) = W - \xi_W \nabla_W \sum_{i=1}^N a_i l(x_i, y_i, W)$$

Approximation of $\widetilde{V}(\widetilde{W}(A))$:

$$\begin{aligned} \widetilde{V}(\widetilde{W}(A)) &\doteq \overline{V}(\overline{W}(A)) \\ &= V - \xi_v \nabla_v \sum_{i=1}^M l(u_i, g(u_i, \overline{W}(A)); V) \end{aligned} \tag{1}$$

For the validation stage:

$$\min_A L(D_{val}, \overline{V}(\overline{W}(A)))$$

$$A = A - \xi_A \nabla_A \overline{V}(\overline{W}(A)) \cdot \nabla_{\overline{V}} L(D_{val}, \overline{V}) \tag{2}$$

$$\begin{aligned} \nabla_A \overline{V}(\overline{W}(A)) &= \nabla_A (V - \xi_v \nabla_v \sum_{i=1}^M l(u_i, g(u_i, \overline{W}(A)); V)) \\ &= \xi_v \nabla_v \nabla_A \sum_{i=1}^M l(u_i, g(u_i, \overline{W}(A)); V) \end{aligned}$$

$$= \xi_V \cdot [\nabla_V \nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V) \cdot \nabla_A \bar{W}(A)]$$

now equation (2) becomes:

$$A = A - \xi_A \xi_V \cdot \nabla_V \nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V) \cdot \nabla_A \bar{W}(A) \cdot \nabla_{\bar{V}} L(D_{val}, \bar{V}) \quad (3)$$

Using Finite difference method to approximate

$$\begin{aligned} & \nabla_V \nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V) \cdot \nabla_{\bar{V}} L(D_{val}, \bar{V}) \\ &= \frac{1}{2\alpha_W} (\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^+) - \nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^-)) \end{aligned}$$

where:

$$\begin{aligned} V^\pm &= V \pm \alpha_W \nabla_{\bar{V}} L(D_{val}, \bar{V}) \\ \alpha_W &= \frac{0.01}{\|\nabla_{\bar{V}} l(D_{val}, \bar{V})\|_2} \end{aligned}$$

Second term in (3):

$$\begin{aligned} \nabla_A \bar{W}(A) &= \nabla_A (W - \xi \nabla_W \sum_{i=1}^N (\alpha_i l(x_i, y; W))) \\ &= \xi_W \nabla_W \nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W) \end{aligned}$$

now equation (3) becomes:

$$A = A - \xi_A \xi_V \left[\frac{1}{2\alpha_W} (\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^+) - \nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^-)) \right] \cdot [\xi_W \nabla_W \nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W)]$$

$$A = A - \xi_W \xi_A \xi_V \frac{1}{2\alpha_W} [(\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^+) - \nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^-))] \cdot [\nabla_W \nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W)]$$

Using Finite difference method to approximate

$$[\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^+)] \cdot [\nabla_W \nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W)]$$

$$= \frac{1}{2\alpha_A} (\nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W^+) - \nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W^-))$$

where:

$$W^\pm = W \pm \alpha_A (\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^+))$$

$$\alpha_A = \frac{0.01}{\|\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^+)\|_2}$$

Using Finite difference method to approximate

$$[\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^-)] \cdot [\nabla_W \nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W)]$$

$$= \frac{1}{2\alpha_A} (\nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W^+) - \nabla_A \sum_{i=1}^n \alpha_i l(x_i, y_i; W^-))$$

where:

$$W^\pm = W \pm \alpha_A (\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^-))$$

$$\alpha_A = \frac{0.01}{\|\nabla_{\bar{W}} \sum_{i=1}^M l(u_i, g(u_i, \bar{W}); V^-)\|_2}$$