

~~All parameters~~

~~the labels~~

$$5a) CE^i = -\sum_{k=1}^K y_k \log \hat{y}_k^i = -\log \hat{y}_e^i = -\log \left[\frac{\hat{e}^{z_e^i}}{\sum_{j=1}^K \hat{e}^{z_j^i}} \right] \quad \text{ignoring superscript of } e$$

$$\frac{\partial CE^i}{\partial z_q^i} = -\frac{1}{\hat{y}_e^i} \frac{\cancel{\hat{e}^{z_q^i}} \delta_{qe} \sum_{j=1}^K \hat{e}^{z_j^i}}{\cancel{\hat{e}^{z_e^i}}} = \frac{+1}{\hat{y}_e^i} \left(\cancel{\delta_{qe} \hat{y}_q^i} + \cancel{\hat{y}_e^i \hat{y}_q^i} \right)$$

$$= \cancel{\delta_{qe}} + \hat{y}_q^i = \underline{-\hat{y}_q^i + \hat{y}_q^i} \Rightarrow \underline{\nabla_{z_q^i} CE^i = -\hat{y}_q^i + \hat{y}_q^i} \in \mathbb{R}^K$$

5b) Need to compute the gradients wrt w^i, w^2, b^i, b^2 .

$$b^i: \frac{\partial CE^i}{\partial b_a^2} = \frac{\partial CE^i}{\partial z_q^i} \frac{\partial z_q^i}{\partial b_a^2} = (\hat{y}_q^i - \hat{y}_a^i) \cdot \delta_{aq} = \underline{\hat{y}_a^i - \hat{y}_a^i}$$

$$w^2: \frac{\partial CE^i}{\partial W_{qp}^2} = \frac{\partial CE^i}{\partial z_q^i} \frac{\partial z_q^i}{\partial W_{qp}^2} = (\hat{y}_q^i - \hat{y}_q^i) \cdot (\delta_{qp} \alpha_p) = \underline{(\hat{y}_q^i - \hat{y}_q^i) \alpha_p^i}$$

$$b^2: \frac{\partial CE^i}{\partial b_a^2} = \frac{\partial CE^i}{\partial z_q^i} \cdot \frac{\partial z_q^i}{\partial a_r^i} \frac{\partial a_r^i}{\partial b_a^2} = (\hat{y}_q^i - \hat{y}_q^i) \cdot W_{qr}^2 \sigma_r^i (1 - \sigma_r^i) \cancel{\delta_{ra}}$$

$$= \underline{(\hat{y}_q^i - \hat{y}_q^i) W_{qr}^2 \sigma_r^i (1 - \sigma_r^i)}$$

$$w^i: \frac{\partial CE^i}{\partial W_{ap}^i} = \frac{\partial CE^i}{\partial z_q^i} \frac{\partial z_q^i}{\partial a_r^i} \frac{\partial a_r^i}{\partial W_{ap}^i} = (\hat{y}_q^i - \hat{y}_a^i) W_{qr}^2 \sigma_r^i (1 - \sigma_r^i) (\delta_{ra} X_\beta^i)$$

$$= \underline{(\hat{y}_q^i - \hat{y}_q^i) W_{qr}^2 \sigma_r^i (1 - \sigma_r^i) X_\beta^i} = (\hat{y}_q^i - \hat{y}_q^i) W_{qr}^2 \alpha_a^i (1 - \alpha_a^i) X_\beta^i$$

~~Final answer~~ Now done - my code works good!

5c) Both regularised & non-regularised converge well.
But the remarkable thing here is that adding the regularization had almost no effect on the training set accuracy & loss, but while giving a significant several percentage point increase in the test dev set accuracy loss. This is due to regularization punishing anything (albeit at only ~~such~~ lower) to a modest increase in bias.

5d) 93% vs 97% for the test set, without and with regularization. This makes sense for the same reasons as above.