

~~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 .

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

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

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

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

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

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

~~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.

$$6a) p(\theta|x,y) = \frac{p(\theta,xy)}{p(xy)} \stackrel{\text{def}}{=} p(y|x,\theta) \frac{p(\theta|y)}{p(y|x)} = p(y|x,\theta) \frac{p(\theta)p(x)}{p(y|x)p(x)}$$

$p(y|x)$ indep. of θ .

$$\downarrow \hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} \, p(y|x,\theta) \frac{p(\theta)}{p(y|x)} = \underset{\theta}{\operatorname{argmax}} \, p(y|x,\theta)p(\theta)$$

$$6b) \text{ If we set our prior to be } p(\theta|\mathcal{O}, \mu^2 I) = \frac{1}{\sqrt{2\pi^k}} \exp\left(-\frac{1}{2\mu^2}\theta^T\theta\right)$$

$$\text{then we get that } \hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} \, p(y|x,\theta) \sim \exp(-\cdot)$$

$$\begin{aligned} \hat{\theta}_{MAP} &= \underset{\theta}{\operatorname{argmax}} \log p(y|x,\theta) - \log(p(y|x,\theta)) - \frac{\theta^T\theta}{2\mu^2} \\ &= \underset{\theta}{\operatorname{argmin}} -\log p(y|x,\theta) + (\frac{1}{2\mu^2})\|\theta\|^2, \quad \lambda = \frac{1}{2\mu^2} \end{aligned}$$

$$6c) \text{ In other words, we have } p(\vec{y}|X,\theta) \sim \exp\left(-\frac{1}{2\sigma^2}(\vec{y}-X\theta)^T(\vec{y}-X\theta)\right)$$

$$\Rightarrow \hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} (\vec{y}-X\theta)^T(\vec{y}-X\theta) + \frac{\theta^T\theta}{2\mu^2} = \underset{\theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \left[\vec{y}^T\vec{y} + \theta^T X^T X \theta - 2\theta^T X^T \vec{y} + \frac{1}{2\mu^2} \theta^T \theta \right]$$

Taking the derivative w.r.t θ of the loss gives :

$$\frac{\partial}{\partial \theta} \text{Loss} = \frac{1}{2} (x^T x \Theta - 2x^T y) + \frac{\Theta}{\sigma^2} = 0$$

$$\Leftrightarrow \cancel{(x^T x \Theta - 2x^T y)} + (\frac{\sigma^2}{2} x^T x + \frac{\lambda}{2} \Theta^T \Theta) \Theta_{MAP} = x^T y \cancel{\Theta}$$

$$\Rightarrow \cancel{\Theta^T \Theta} = \cancel{x^T x \Theta - 2x^T y} + \cancel{\Theta^T x^T y} \quad \Theta_{MAP} = \underline{\left[x^T x + \frac{\sigma^2}{2} \mathbb{I} \right]^{-1}} \underline{x^T y}$$

$$= \underline{\left(x^T x + 2\sigma^2 \lambda \mathbb{I} \right) x^T y}$$

6d) Now $\underline{\Theta_{MAP}} = \underset{\Theta}{\operatorname{argmin}} -\log p(y|x,\Theta) + \sum_{i=1}^k \frac{1}{b} |\Theta_i|_1$

$$= \underset{\Theta}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\vec{y} - X\vec{\Theta}\|_2^2 + \frac{1}{b} |\Theta|_1$$

$$= \underset{\Theta}{\operatorname{argmin}} \|\vec{y} - X\vec{\Theta}\|_2^2 + \frac{\alpha\sigma^2}{b} |\Theta|_1, \text{ with } \gamma = \frac{2\sigma^2}{b}$$