

1. (Math) Recall that when ℓ_1 regularization is used, and assuming that the Hessian matrix is diagonal with positive entries, the objective function can be approximated by

$$\hat{L}_R(\theta) := \sum_{i=1}^d \left[\frac{1}{2} H_{ii} (\theta_i - \theta_i^*)^2 + \alpha |\theta_i| \right]$$

Solve this in the close form expression: show that the optimal solution θ_R^* for the objective $\hat{L}_R(\theta)$ is that as shown on slide 26 of "Deep learning lecture 3: Regularization I".

Hint: note that $\alpha/H_{ii} > 0$.

$$\text{Proof: } \hat{L}_R(\theta_R^*) = \sum_{i=1}^d \left[\frac{1}{2} H_{ii} (\theta_{Ri} - \theta_{Ri}^*)^2 + \alpha |\theta_{Ri}| \right] = \hat{L}(\theta_R^*) + \alpha \|\theta_R^*\|_1$$

$$\begin{aligned} (\theta_R^*)_i &= \frac{d \hat{L}_R}{d \theta} = \frac{d}{d \theta} \left[\frac{1}{2} H_{ii} (\theta_i - \theta_i^*)^2 \right] + \frac{d}{d \theta} [\alpha |\theta_i|] \\ &= \theta_i^* - \frac{\alpha}{H_{ii}} \end{aligned}$$

$$\text{Since } \theta_R^* \text{ is optimal, } \theta_R^* = \theta^* - \frac{\alpha}{H_{ii}} = 0, \quad \theta_R^* = \theta^* = \frac{\alpha}{H_{ii}}$$

$$\hat{L}(\theta) = L(\theta^*) + \frac{1}{2} (\theta - \theta^*)^T H (\theta - \theta^*)$$

$\therefore H$ matrix is positive and diagonal

$$\begin{bmatrix} +\frac{d^2}{dx^2} & 0 \\ 0 & +\frac{d^2}{dy^2} \end{bmatrix}$$

$$\therefore (\theta - \theta^*)^T H (\theta - \theta^*) = H (\theta - \theta^*)^2$$

$$\hat{L}(\theta) = L(\theta^*) + \frac{1}{2} H (\theta - \theta^*)^2$$

By drop the constant $L(\theta^*)$,

$$\hat{L}(\theta) = \frac{1}{2} H (\theta - \theta^*)^2$$

$$\sum \alpha |\theta_i| = \alpha \sum |\theta_i| = \alpha \|\theta\|_1$$

$$\begin{aligned} \hat{L}_R(\theta_R^*) &= \sum_{i=1}^d \left[\frac{1}{2} H_{ii} (\theta_{Ri} - \theta_{Ri}^*)^2 + \alpha |\theta_{Ri}| \right] \\ &= \hat{L}(\theta_R^*) + \alpha \|\theta_R^*\|_1 \end{aligned}$$

= slide 26

2. (Math) Consider a three layer network:

$$h^1 = \sigma(W^1 x), \quad h^2 = \sigma(W^2 h^1), \quad f(x) = \langle w^3, h^2 \rangle.$$

See Figure 1 for an illustration. Compute $\frac{\partial f}{\partial w_{ij}^1}$

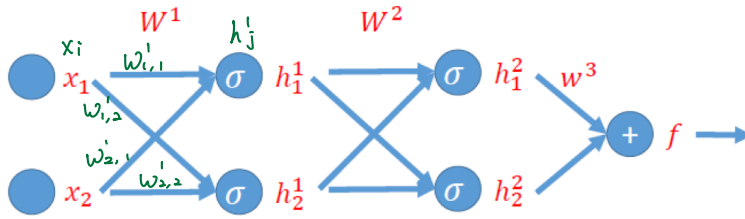


Figure 1: An illustration of the three layer network

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma'(x) = \frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x))$$

$$\begin{aligned} \frac{\partial f}{\partial w_{ij}^1} &= \frac{\partial f}{\partial h_1^2} \cdot \frac{\partial h_1^2}{\partial \text{net}_1^2} \cdot \frac{\partial \text{net}_1^2}{\partial h_j^1} \cdot \frac{\partial h_j^1}{\partial \text{net}_j^1} \cdot \frac{\partial \text{net}_j^1}{\partial w_{ij}^1} \\ &+ \frac{\partial f}{\partial h_2^2} \cdot \frac{\partial h_2^2}{\partial \text{net}_2^2} \cdot \frac{\partial \text{net}_2^2}{\partial h_j^1} \cdot \frac{\partial h_j^1}{\partial \text{net}_j^1} \cdot \frac{\partial \text{net}_j^1}{\partial w_{ij}^1} \\ &= w^3 \times \sigma'(\text{net}_1^2) \times w^2 \times \sigma'(\text{net}_j^1) \times x_i \\ &+ w^3 \times \sigma'(\text{net}_2^2) \times w^2 \times \sigma'(\text{net}_j^1) \times x_i \\ &= [\sigma'(\text{net}_1^2) + \sigma'(\text{net}_2^2)] \times w^3 w^2 \sigma'(\text{net}_j^1) x_i \end{aligned}$$

$$\begin{aligned} \sigma'(\text{net}_1^2) &= \sigma(\text{net}_1^2)(1 - \sigma(\text{net}_1^2)) \\ &= h_1^2(1 - h_1^2) \\ \sigma'(\text{net}_2^2) &= h_2^2(1 - h_2^2) \\ \sigma'(\text{net}_j^1) &= h_j^1(1 - h_j^1) \end{aligned}$$

$$= [h_1^2(1 - h_1^2) + h_2^2(1 - h_2^2)] \cdot w^3 w^2 \cdot h_j^1(1 - h_j^1) \cdot x_i$$