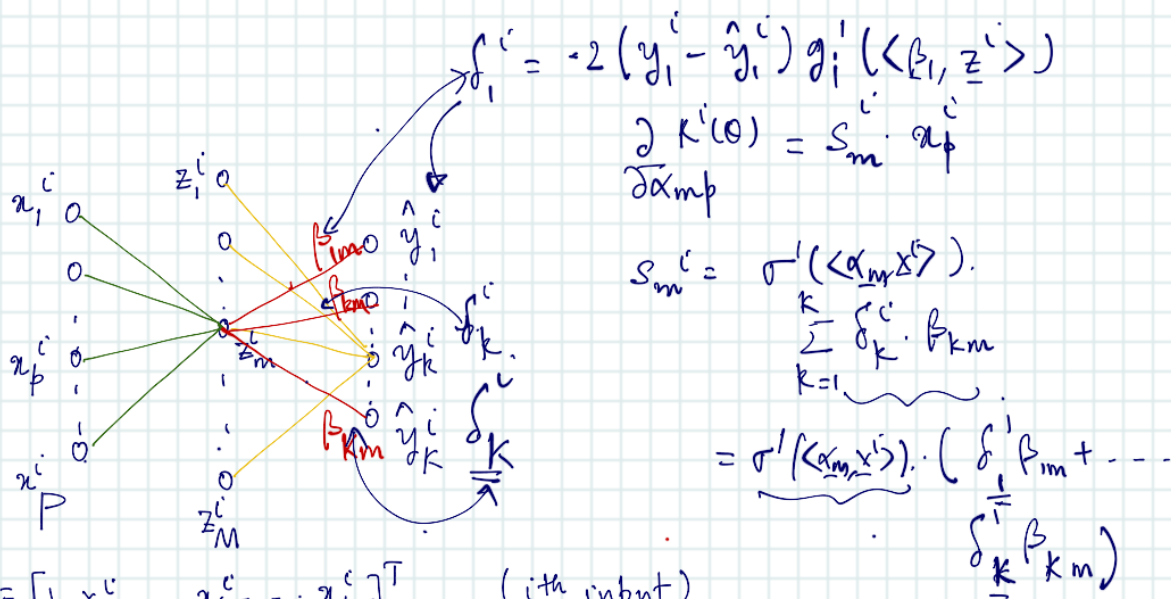


14/9/23

Deep Learning

- Recap
- Minimizing ANN loss
- Back propagation



- Input vector: $\underline{x}^i = [1, x_1^i, \dots, x_p^i]^\top$ (i^{th} input)
- Hidden layer vector: $\underline{z}^i = [1, z_1^i, \dots, z_m^i, \dots, z_n^i]^\top$
- Output vector: $\hat{\underline{y}}^i = [\hat{y}_1^i, \dots, \hat{y}_k^i, \dots, \hat{y}_K^i]^\top$
- First layer weights: $\underline{\alpha}_m = [\alpha_{m0}, \alpha_{m1}, \dots, \alpha_{mp}]^\top$ (m^{th} hidden layer node)
- Second layer weights: $\underline{\beta}_k = [\beta_{k0}, \beta_{k1}, \dots, \beta_{km}]^\top$ (k^{th} output layer node).
- Output at m^{th} hidden layer node: $z_m^i = \sigma(\langle \underline{\alpha}_m, \underline{x}^i \rangle)$
- Output at k^{th} output layer node: $\hat{y}_k^i = g_k(\langle \underline{\beta}_k, \underline{z}^i \rangle)$
- Risk (or loss): $R(\theta) = \sum_{i=1}^n d(y^i, \hat{y}^i) = \sum_{i=1}^n R^i(\theta)$

$$\theta = \{ \underline{\alpha}_1, \dots, \underline{\alpha}_m, \underline{\beta}_1, \dots, \underline{\beta}_K \}$$

$$= \sum_{i=1}^n \sum_{k=1}^K d(y_k^i, \hat{y}_k^i)$$

$$\text{MSE loss: } R(\theta) = \sum_{i=1}^n \sum_{k=1}^K (y_k^i - \hat{y}_k^i)^2 = \sum_{i=1}^n R^i(\theta) \text{ where } R^i(\theta) = \sum_{k=1}^K (y_k^i - \hat{y}_k^i)^2$$

- Find θ that gives us a local optimum

$$\theta^{(r)} = \theta^{(r-1)} - \eta^{(r)} \cdot \nabla_{\theta^{(r-1)}} R(\theta)$$

$$\frac{\partial R^i(\theta)}{\partial \beta_{km}} = -2(y_k^i - \hat{y}_k^i) g_k'(\langle \beta_k, \underline{z}^i \rangle) \cdot z_m^i$$

$$J(\theta) = \sum_{i=1}^n R^i(\theta)$$

$$\frac{\partial R(\theta)}{\partial \beta_{km}} = \sum_{i=1}^n \frac{\partial R^i(\theta)}{\partial \beta_{km}} = -2 \sum_{i=1}^n (y_k^i - \hat{y}_k^i) \cdot g'_k(\langle \beta_k, z^i \rangle) \cdot z_m^i$$

• $\frac{\partial R^i(\theta)}{\partial \beta_{km}} = \delta_k^i z_m^i$, where $\delta_k^i = -2 (y_k^i - \hat{y}_k^i) g'_k(\langle \beta_k, z^i \rangle)$ — (1)

• $\frac{\partial R^i(\theta)}{\partial \alpha_{mp}} = -2 \sum_{k=1}^K \frac{(y_k^i - \hat{y}_k^i) \cdot g'_k(\langle \beta_k, z^i \rangle) \cdot \beta_{km} \cdot \sigma'(\langle \alpha_m, x^i \rangle) \cdot x_p^i}{\beta_{km} \cdot \sigma'(\langle \alpha_m, x^i \rangle)}$

$$\frac{\partial R^i(\theta)}{\partial \alpha_{mp}} = s_m^i x_p^i ; s_m^i = \sigma'(\langle \alpha_m, x^i \rangle) \cdot \sum_{k=1}^K \delta_k^i \beta_{km}$$

• Backpropagation: The "error" at the output node is backpropagated to find the partial derivative $\frac{\partial R^i(\theta)}{\partial \alpha_{mp}}$, and thereby the gradient

$$\theta^{(r)} = \theta^{(r-1)} - \eta^{(r)} \nabla_{\theta^{(r-1)}} R(\theta)$$

$$\boxed{\theta^{(r)} = \theta^{(r-1)} - \underbrace{\eta^{(r)}}_{\text{learning rate}} \nabla_{\theta^{(r-1)}} \left(\sum_{i=1}^n R^i(\theta) \right)} \quad - (2)$$

• θ is randomly initialized to small values.