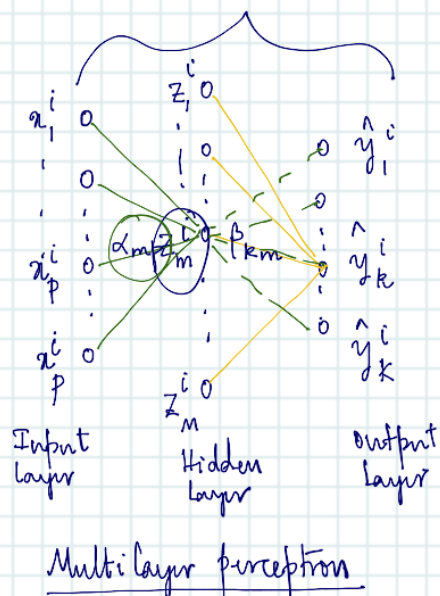


- Recap ✓
- Loss function ✓
- Gradient descent ✓
- Back propagation

• Recap:



data point ground truth
↓ ↓
 $\mathcal{D} = \{(\underline{x}^i, \underline{y}^i), \dots, (\underline{x}^n, \underline{y}^n)\}$

$$\underline{x}^i = [1 \ x_1^i \ \dots \ x_p^i]^T$$

$$\underline{z}^i = [1 \ z_1^i \ \dots \ z_m^i]^T$$

$$\underline{\hat{y}}^i = [\hat{y}_1^i \ \dots \ \hat{y}_k^i]^T$$

$$\underline{\alpha}_m = [\alpha_{m0}, \alpha_{m1} \ \dots \ \alpha_{mp}]^T$$

↑
bias

$$\underline{\beta}_k = [\beta_{k0}, \beta_{k1} \ \dots \ \beta_{km}]^T$$

$$z_m^i = \sigma(\langle \underline{\alpha}_m, \underline{x}^i \rangle)$$

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

$$\hat{y}_k^i = \text{softmax}(\langle \underline{\beta}_k, \underline{z}^i \rangle) = g_k(\langle \underline{\beta}_k, \underline{z}^i \rangle)$$

• Loss function: $R(\theta) \leftrightarrow \mathcal{L}(\theta) = \sum_{i=1}^n d(\underline{y}^i, \underline{\hat{y}}^i)$

↓ ↓
Risk Loss

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

For simplicity, let $d(y_k^i, \hat{y}_k^i) = (y_k^i - \hat{y}_k^i)^2$

$$\therefore R(\theta) = \sum_{i=1}^n \sum_{k=1}^K (y_k^i - \hat{y}_k^i)^2$$

$$R(\theta) = \sum_{i=1}^n \sum_{k=1}^K (y_k^i - \underbrace{f_k(\underline{x}^i; \theta)}_{\text{machine}})^2$$

- Find $\underline{\theta}$ that minimizes the loss or risk $R(\theta)$; $\theta^* = \underset{\theta \in \Theta}{\operatorname{argmin}} R(\theta)$
- Exercise: Show that $f_k(x^i; \theta)$ is a non-convex function of θ .
- We will rely on gradient based techniques for find a local optimum.

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

$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_m \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix}^{(r)} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_m \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}^{(r-1)} - \eta^{(r)} \begin{bmatrix} \frac{\partial R(\theta)}{\partial \alpha_1} \\ \vdots \\ \frac{\partial R(\theta)}{\partial \alpha_m} \\ \frac{\partial R(\theta)}{\partial \beta_{km}} \\ \frac{\partial R(\theta)}{\partial \beta_{km}} \\ \frac{\partial R(\theta)}{\partial \beta_{km}} \end{bmatrix}^{(r-1)}$$

$$\text{Find } \frac{\partial R(\theta)}{\partial \beta_{km}} = \frac{\partial}{\partial \beta_{km}} \sum_{i=1}^n \sum_{j=1}^k (y_j^i - \hat{y}_j^i)^2$$

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

↳ what happens when g_k is softmax?

$$\frac{\partial R(\theta)}{\partial \alpha_{mp}} = \frac{\partial}{\partial \alpha_{mp}} \sum_{i=1}^n \sum_{k=1}^k (y_k^i - \hat{y}_k^i)^2$$

$$= -2 \sum_{i=1}^n \sum_{k=1}^k (y_k^i - \hat{y}_k^i) g_k'(\langle \underline{\beta}_k, \underline{z}^i \rangle) \cdot \beta_{km} \sigma'(\langle \underline{\alpha}_m, \underline{x}^i \rangle) \cdot x_p^i$$

Heart of back prop