

6. K-class classification \rightarrow output space $\mathcal{Y} = \{1, 2, \dots, K\}$

Matrix W represents a hypothesis $h_W(\cdot)$

$$W = [\vec{w}_1 \vec{w}_2 \dots \vec{w}_k \dots \vec{w}_K]_{d \times K} \quad \vec{w}_k \in \mathbb{R}^d \quad \forall k \in \{1, \dots, K\}$$

$$h_W(\vec{x}) = \frac{e^{\vec{w}_y^T \vec{x}}}{\sum_{k=1}^K e^{\vec{w}_k^T \vec{x}}} \xrightarrow{\text{approx.}} p(y|\vec{x})$$

$\mathcal{D} = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2), \dots, (\vec{x}_N, y_N)\} \stackrel{iid}{\sim} p(\vec{x})$, target distribution $p(y|\vec{x})$

$$\mathcal{L}(h_W | \mathcal{D}) = p(\mathcal{D} | h_W) \propto \prod_{n=1}^N h_{y_n}(\vec{x}_n)$$

$$\begin{aligned} \min -\ln \text{likelihood} &\propto \min -\ln \prod_{n=1}^N h_{y_n}(\vec{x}_n) = \min \sum_{n=1}^N -\ln h_{y_n}(\vec{x}_n) = \min \sum_{n=1}^N \text{ew}(\vec{w}, \vec{x}_n, y_n) \\ &= \min \sum_{n=1}^N -\ln \frac{e^{\vec{w}_{y_n}^T \vec{x}_n}}{\sum_{k=1}^K e^{\vec{w}_k^T \vec{x}_n}} \end{aligned}$$

$$\frac{\partial \text{ew}(\vec{w}, \vec{x}, y)}{\partial w_i} = \frac{\partial}{\partial w_i} \left[-\vec{w}_y^T \vec{x} + \ln \left(\sum_{k=1}^K e^{\vec{w}_k^T \vec{x}} \right) \right]$$

$$= \frac{\partial}{\partial w_i} (-\vec{w}_y^T \vec{x}) + \frac{\partial}{\partial w_i} \left(\sum_{k=1}^K e^{\vec{w}_k^T \vec{x}} \right)$$

$$= \begin{cases} -\vec{x} + h_i(\vec{x}) \cdot \vec{x} & y=i \\ h_i(\vec{x}) \cdot \vec{x} & y \neq i \end{cases}$$

$$\left(\frac{\partial \text{ew}(\vec{w}, \vec{x}_n, y_n)}{\partial w_i} = \begin{cases} -\vec{x}_n + h_i(\vec{x}_n) \cdot \vec{x}_n & y_n=i \\ h_i(\vec{x}_n) \cdot \vec{x}_n & y_n \neq i \end{cases} \right)$$

$$\frac{\partial}{\partial w_i} (-\vec{w}_y^T \vec{x}) \rightarrow \text{if } y=i, \text{ then } \frac{\partial}{\partial w_i} (-\vec{w}_y^T \vec{x}) = \frac{\partial}{\partial w_i} (-\vec{w}_i^T \vec{x}) = -\vec{x}$$

$$\text{if } y \neq i, \text{ then } \frac{\partial}{\partial w_i} (-\vec{w}_y^T \vec{x}) = 0$$

$$\frac{\partial}{\partial w_i} \left(\sum_{k=1}^K e^{\vec{w}_k^T \vec{x}} \right) \stackrel{\text{chain rule}}{=} \frac{1}{\sum_{k=1}^K e^{\vec{w}_k^T \vec{x}}} \cdot e^{\vec{w}_i^T \vec{x}} \cdot \vec{x} = h_i(\vec{x}) \cdot \vec{x}$$

Let $V = [\vec{v}_1 \vec{v}_2 \dots \vec{v}_K]_{d \times K}$ where $\vec{v}_i \in \mathbb{R}^d \quad \forall i \in \{1, \dots, K\}$

\vec{v}_1 \rightarrow update direction for \vec{w}_1

$$\forall i \in \{1, \dots, K\} \quad i \neq y_n, \quad \vec{v}_i = -h_i(\vec{x}_n) \cdot \vec{x}_n$$

$$\text{for } i = y_n, \quad \vec{v}_i = \vec{v}_{y_n} = -[-\vec{x}_n + h_i(\vec{x}_n) \vec{x}_n] = \vec{x}_n - h_i(\vec{x}_n) \vec{x}_n = \vec{x}_n (1 - h_i(\vec{x}_n))$$

$$\text{i.e. } V = [\vec{x}_n - h_1(\vec{x}_n) \vec{x}_n \dots \vec{x}_n (1 - h_i(\vec{x}_n)) \dots \vec{x}_n - h_K(\vec{x}_n) \vec{x}_n]$$

$$= \vec{x}_n [-h_1(\vec{x}_n) \dots -h_i(\vec{x}_n) (1 - h_i(\vec{x}_n)) -h_i(\vec{x}_n) \dots -h_K(\vec{x}_n)]$$

$$= \vec{x}_n \cdot \vec{u}^T$$

$$\therefore \vec{u}^T = [-h_1(\vec{x}_n) \dots (1 - h_i(\vec{x}_n)) \dots -h_K(\vec{x}_n)]$$

$$\Rightarrow \vec{u} = [-h_1(\vec{x}_n) \dots (1 - h_i(\vec{x}_n)) \dots -h_K(\vec{x}_n)]^T \in \mathbb{R}^{K \times 1} \quad \square$$