

For one data point  $i$ , its data loss is computed as follow:

$$L_i = \sum_{j \neq y_i}^{C-1} \max(0, s_j - s_{y_i} + 1) \quad \text{where } s_j = X_i W_j$$

score for class  $j$

$$s_j: 1 \times 1 = X_i: 1 \times D * W_j: D \times 1$$

For its scores for all classes,

$$S[i]: 1 \times C = X_i: 1 \times D * W: D \times C$$

For all data points,

$$S: \begin{matrix} C \\ \text{num of train data} \end{matrix} = X: \begin{matrix} D \\ \text{num of train data} \end{matrix} * W: \begin{matrix} C \\ D \end{matrix}$$

To compute gradient

$$\nabla_{s_j} L_i = \mathbb{1}(s_j - s_{y_i} + 1 > 0) \text{ where } \mathbb{1}(\text{true}) = 1$$

$$\mathbb{1}(\text{false}) = 0$$

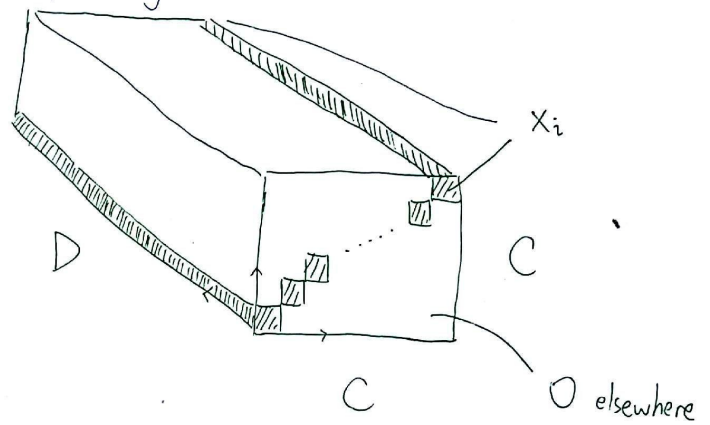
$$\nabla_{s_{y_i}} L_i = - \sum_{j \neq y_i} \mathbb{1}(s_j - s_{y_i} + 1 > 0)$$

$$\nabla_{s_{[i]}} L_i = \frac{C}{1}$$

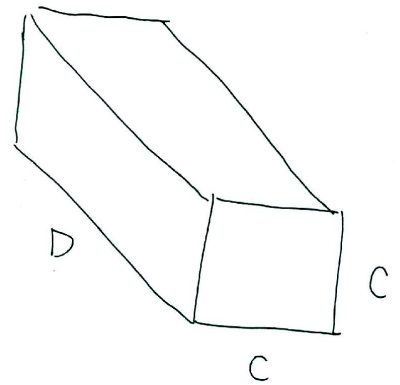
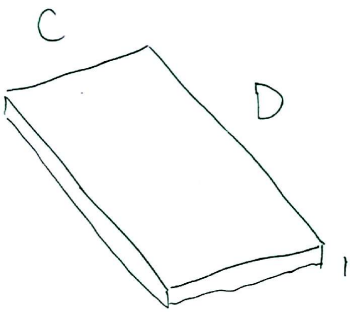
$$\nabla_{w_j} S_j = x_i$$

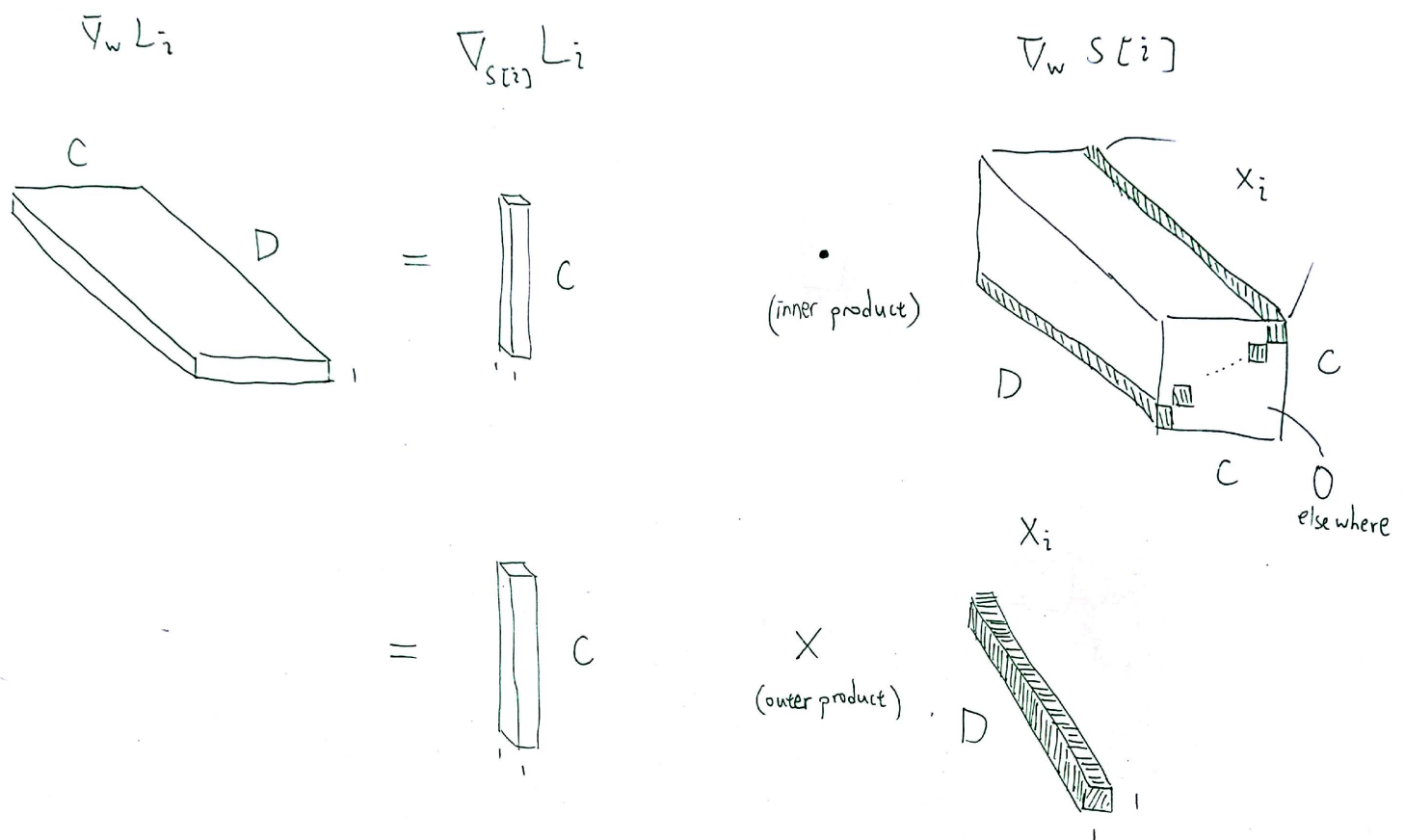
$$\nabla_{w_j} S_k = 0 \quad k \neq j$$

$$\nabla_w S[i] =$$



$$\nabla_w L_i = \nabla_{S(i)} L_i * \nabla_w S[i]$$





... which simplifies computation a bit.