

$$x \in \mathbb{R}^{n \times 1} \quad W \in \mathbb{R}^{m \times n} \quad E \in \mathbb{R}^{m \times d}$$

$$s = Wx \in \mathbb{R}^{m \times 1}$$

$$\hat{s} = \text{argtopk}(s) \Rightarrow s \text{ 中 topk 元素的 index } \in \mathbb{R}^{k \times 1}, \text{ 其每个元素 } \in [0, \dots, m-1]$$

$$\text{hard: } \text{res} = \underbrace{\text{one-hot}(\hat{s})}_\times E = E[\hat{s}] \in \mathbb{R}^{k \times d}$$

$$\text{loss} = (\text{res} \times \times 2). \text{sum}()$$

$\text{loss.backward}() \Rightarrow$  对  $E$  有梯度 (仅从 topk 对应的 entry), 对  $W$  没梯度

$$\text{soft: } \text{res} = \text{one-hot}(\hat{s}) E \odot \underbrace{\sigma(s)}_{\begin{cases} 0 & s_i < \gamma \\ s_i & s_i \geq \gamma \end{cases}}$$

$$\text{loss} = (\text{res} \times \times 2). \text{sum}()$$

$\text{loss.backward}() \Rightarrow$  对  $E$  有梯度 (同上), 对  $W$  有梯度 (仅从  $s_i \geq \gamma$  的 entry)