

Attention (self/cross, hard/soft)

Dealing with sets

$$\mathbf{h} = \mathbf{X} \mathbf{a}$$

Self-attention (I)

$$\{\mathbf{x}_i\}_{i=1}^t = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t\} \rightsquigarrow \mathbf{X} \in \mathbb{R}^{n \times t}, \quad \mathbf{x}_i \in \mathbb{R}^n$$

$$\mathbf{h} = \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \dots + \alpha_t \mathbf{x}_t = \mathbf{X} \mathbf{a} \in \mathbb{R}^n$$

$$\alpha_i > 0$$

$$\mathbf{X} \doteq \begin{bmatrix} | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_t \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times t}$$

soft attention: $\mathbf{a}^\top \mathbf{1} = 1$

hard attention: $\mathbf{a} \in \mathbb{I}_t$

$$\mathbf{h} = \mathbf{X} \mathbf{a}$$

Self-attention (II)

$$\mathbf{a} = \text{softmax}_{\beta}(\mathbf{X}^{\top} \mathbf{x}) \in \mathbb{R}^t$$

$$\{\mathbf{x}_i\}_{i=1}^t \rightsquigarrow \{\mathbf{a}_i\}_{i=1}^t \rightsquigarrow \mathbf{A} \in \mathbb{R}^{t \times t}$$

$$\{\mathbf{a}_i\}_{i=1}^t \rightsquigarrow \{\mathbf{h}_i\}_{i=1}^t \rightsquigarrow \mathbf{H} \in \mathbb{R}^{n \times t}$$

$$\mathbf{H} = \mathbf{X} \mathbf{A} \in \mathbb{R}^{n \times t}$$

Key-value store

- Paradigm for
 - storing (saving)
 - retrieving (querying)
 - managing
- an associative array (dictionary / hash table)

Queries, keys, and values $\{\mathbf{q}_i\}_{i=1}^t \rightsquigarrow \mathbf{Q} \in \mathbb{R}^{d' \times t}$

$$\mathbf{q} = \mathbf{W}_{\mathbf{q}} \mathbf{x}, \quad \mathbf{k} = \mathbf{W}_{\mathbf{k}} \mathbf{x}, \quad \mathbf{v} = \mathbf{W}_{\mathbf{v}} \mathbf{x} \quad \beta = \frac{1}{\sqrt{d'}}$$

$$\mathbf{q}, \mathbf{k} \in \mathbb{R}^{d'}, \quad \mathbf{v} \in \mathbb{R}^{d''}$$

$$\{\mathbf{x}_i\}_{i=1}^t \rightsquigarrow \{\mathbf{q}_i\}_{i=1}^t, \{\mathbf{k}_i\}_{i=1}^t, \{\mathbf{v}_i\}_{i=1}^t \rightsquigarrow \mathbf{Q}, \mathbf{K}, \mathbf{V}$$

$$\mathbf{a} = \text{softargmax}_{\beta}(\mathbf{K}^{\top} \mathbf{q}) \in \mathbb{R}^t \quad \mathbf{h} = \mathbf{V} \mathbf{a} \in \mathbb{R}^{d''}$$

$$\{\mathbf{q}_i\}_{i=1}^t \rightsquigarrow \{\mathbf{a}_i\}_{i=1}^t \rightsquigarrow \mathbf{A} \in \mathbb{R}^{t \times t} \quad \mathbf{H} = \mathbf{V} \mathbf{A} \in \mathbb{R}^{d'' \times t}$$