

# Attention via $\log \sum \exp$ energy

Alexander Tschantz

January 20, 2025

## 1 General Framework

We consider a directed graph of  $A$  nodes, where each node is a vector  $\{\mathbf{v}_a\}_{a=1}^A$  with  $\mathbf{v}_a \in \mathbb{R}^d$ . Each node  $\mathbf{v}_a$  has a set of parents  $\mathcal{P}(a) \subseteq \{1, 2, \dots, A\}$ .

We define a generic *similarity* function

$$\text{sim}(\mathbf{v}_a, \mathbf{v}_p),$$

which measures how well  $\mathbf{v}_a$  is “explained by” or “aligned with” its parent  $\mathbf{v}_p$ . This function may have parameters, as in a dot-product model  $\mathbf{v}_a^\top \mathbf{v}_p$ , or a conditional-probability-based model such as a Gaussian log-likelihood term  $-\frac{1}{2}(\mathbf{v}_a - \mathbf{v}_p)^\top \Sigma^{-1}(\mathbf{v}_a - \mathbf{v}_p)$ .

### 1.1 Energy Function

Each node  $\mathbf{v}_a$  has a  $\log \sum \exp$ -type term over its parents. We sum these across all nodes to define the total energy:

$$E(\{\mathbf{v}_a\}) = - \sum_{a=1}^A \ln \left( \sum_{p \in \mathcal{P}(a)} \exp(\text{sim}(\mathbf{v}_a, \mathbf{v}_p)) \right).$$

Informally, each  $\mathbf{v}_a$  seeks to place high mass on those parents  $\mathbf{v}_p$  yielding large similarity scores.

### 1.2 Gradient Updates

For a single node  $\mathbf{v}_a$ , the gradient of the energy decomposes into two parts, reflecting two ways in which  $\mathbf{v}_a$  can appear in the summations:

$$-\frac{\partial E}{\partial \mathbf{v}_a} = \underbrace{\sum_{p \in \mathcal{P}(a)} \text{softmax}_p(\text{sim}(\mathbf{v}_a, \mathbf{v}_p)) \frac{\partial}{\partial \mathbf{v}_a} \text{sim}(\mathbf{v}_a, \mathbf{v}_p)}_{\text{("being explained by its parents")}} + \underbrace{\sum_{c: a \in \mathcal{P}(c)} \text{softmax}_a(\text{sim}(\mathbf{v}_c, \mathbf{v}_a)) \frac{\partial}{\partial \mathbf{v}_a} \text{sim}(\mathbf{v}_c, \mathbf{v}_a)}_{\text{("explaining its children")}}.$$

Above, the summation  $\sum_{p \in \mathcal{P}(a)}$  iterates over the parents of  $a$ , while  $\sum_{c: a \in \mathcal{P}(c)}$  iterates over all children  $c$  such that  $a$  is in their parent set. The softmax is taken over the appropriate parent indices in each case.

### 1.3 Proof of Gradient (Appendix)

A full derivation, with explicit sums over  $\mathcal{P}(c)$ , is given in Appendix A. There we show how collecting terms in the derivative leads precisely to the two-term decomposition above.

## 2 Gaussian Mixture Models (GMMs)

**Setup.** Let  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$ . We consider  $K$  mixture components, each with mean  $\boldsymbol{\mu}_k \in \mathbb{R}^d$  and covariance  $\Sigma_k$ . Defining  $\pi_k$  as the mixing proportion, a standard GMM log-likelihood term can be written in a form that matches our framework.

In particular, define:

$$\mathbf{A}_{ik} = \ln \pi_k - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k) \in \mathbb{R}^{N \times K}.$$

**Energy.**

$$E^{\text{GMM}}(\mathbf{X}, \{\boldsymbol{\mu}_k\}) = - \sum_{i=1}^N \ln \left( \sum_{k=1}^K \exp(\mathbf{A}_{ik}) \right).$$

**Gradient.** If we differentiate w.r.t.  $\boldsymbol{\mu}_k$ , then

$$-\frac{\partial E^{\text{GMM}}}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \text{softmax}_k(\mathbf{A}_{ik}) \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_i - \boldsymbol{\mu}_k).$$

Setting this gradient to zero yields the usual GMM M-step:

$$\boldsymbol{\mu}_k = \frac{\sum_{i=1}^N \text{softmax}_k(\mathbf{A}_{ik}) \mathbf{x}_i}{\sum_{i=1}^N \text{softmax}_k(\mathbf{A}_{ik})}.$$

### 3 Self-Attention

**Setup.** Let  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ , where each  $\mathbf{x}_i \in \mathbb{R}^d$ . We define two learnable weight matrices,  $\mathbf{W}^Q, \mathbf{W}^K \in \mathbb{R}^{d \times d}$ , and construct

$$\mathbf{Q} = \mathbf{W}^Q \mathbf{X}, \quad \mathbf{K} = \mathbf{W}^K \mathbf{X}.$$

Let  $\mathbf{q}_c$  be the  $c$ -th column of  $\mathbf{Q}$  and  $\mathbf{k}_b$  the  $b$ -th column of  $\mathbf{K}$ . Then we define

$$\text{sim}(\mathbf{x}_b, \mathbf{x}_c) \hat{=} \mathbf{k}_b^\top \mathbf{q}_c.$$

On a single line, we gather these into the matrix:

$$\mathbf{A}_{bc} = \mathbf{k}_b^\top \mathbf{q}_c \in \mathbb{R}^{N \times N}.$$

**Energy.**

$$E^{\text{SA}}(\mathbf{X}) = - \sum_{c=1}^N \ln \left( \sum_{b=1}^N \exp(\mathbf{A}_{bc}) \right).$$

**Gradient.** Differentiating w.r.t.  $\mathbf{x}_i$  gives two terms, as each token  $\mathbf{x}_i$  acts both as a “query” for some other tokens and a “key” for yet others:

$$-\frac{\partial E^{\text{SA}}}{\partial \mathbf{x}_i} = \underbrace{\sum_{b=1}^N \text{softmax}_b(\mathbf{A}_{b,i}) \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_b}_{\text{(query side)}} + \underbrace{\sum_{c=1}^N \text{softmax}_i(\mathbf{A}_{i,c}) \mathbf{W}_K^\top \mathbf{W}_Q \mathbf{x}_c}_{\text{(key side)}}.$$

### 4 Cross Attention

**Setup.** In cross attention, we have one set of *query* vectors and a separate set of *key* vectors. Let  $\mathbf{Q} = \mathbf{W}^Q \mathbf{X}^Q \in \mathbb{R}^{d \times N_Q}$  and  $\mathbf{K} = \mathbf{W}^K \mathbf{X}^K \in \mathbb{R}^{d \times N_K}$ , where  $\mathbf{X}^Q$  has  $N_Q$  query tokens and  $\mathbf{X}^K$  has  $N_K$  key tokens. Denote  $\mathbf{q}_c$  as the  $c$ -th query column of  $\mathbf{Q}$  and  $\mathbf{k}_b$  as the  $b$ -th key column of  $\mathbf{K}$ .

We define

$$\mathbf{A}_{b,c} = \mathbf{k}_b^\top \mathbf{q}_c \in \mathbb{R}^{N_K \times N_Q}.$$

This is a matrix of pairwise similarities between keys and queries.

**Energy.**

$$E^{\text{Cross}}(\mathbf{X}^Q, \mathbf{X}^K) = - \sum_{c=1}^{N_Q} \ln \left( \sum_{b=1}^{N_K} \exp(\mathbf{A}_{b,c}) \right).$$

Minimizing this encourages each query  $\mathbf{q}_c$  to place large mass on keys  $\mathbf{k}_b$  that yield higher dot-products.

**Gradient.** As in the self-attention derivation, taking the derivative w.r.t. a single query or key token yields sums weighted by the appropriate softmax terms. For example, w.r.t. the query-side vector  $\mathbf{x}_i^Q$ ,

$$-\frac{\partial E^{\text{Cross}}}{\partial \mathbf{x}_i^Q} = \sum_{b=1}^{N_K} \text{softmax}_b(\mathbf{A}_{b,i}) \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_b^K + \dots$$

and similarly a key  $\mathbf{x}_j^K$  appears in the “explaining children” part of the gradient.

## 5 Hopfield Networks (Softmax Version)

**Setup.** We have data vectors  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ , each  $\mathbf{x}_i \in \mathbb{R}^d$ , and memory vectors  $\mathbf{m}_\mu \in \mathbb{R}^d$  for  $\mu = 1, \dots, K$ . Define

$$\mathbf{A}_{i\mu} = \mathbf{x}_i^\top \mathbf{m}_\mu \in \mathbb{R}^{N \times K}.$$

**Energy.**

$$E^{\text{Hopfield}}(\mathbf{X}) = - \sum_{i=1}^N \ln \left( \sum_{\mu=1}^K \exp(\mathbf{A}_{i\mu}) \right).$$

**Gradient.**

$$-\frac{\partial E^{\text{Hopfield}}}{\partial \mathbf{x}_i} = \sum_{\mu=1}^K \text{softmax}_\mu(\mathbf{A}_{i\mu}) \mathbf{m}_\mu.$$

Hence each  $\mathbf{x}_i$  is updated toward a softmax-weighted combination of the memory vectors.

## 6 Slot Attention

Slot Attention can be seen as cross attention in which we normalize across slots (queries) for each token (key), rather than the usual normalization across the token dimension.

**Setup.** Let  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$  be the set of tokens (e.g. image patches), where each  $\mathbf{x}_j \in \mathbb{R}^d$ . We also have a set of  $S$  latent “slots”:  $\boldsymbol{\mu}_i \in \mathbb{R}^d$  for  $i = 1, \dots, S$ . As in cross attention, we define learnable transforms:

$$\mathbf{W}_K, \mathbf{W}_Q \in \mathbb{R}^{d \times d}.$$

The tokens serve as keys (and potentially values), while the slots serve as queries. However, the key difference is that *each token decides its distribution over slots* (hence the normalization is over  $i$  for each fixed token  $j$ ).

**Energy.** We write the negative log-likelihood as a sum over tokens  $j = 1, \dots, N$ , and in each term we do a  $\log \sum \exp$  over the slots  $i = 1, \dots, S$ . Concretely,

$$E^{\text{Slot}}(\{\boldsymbol{\mu}_i\}) = - \sum_{j=1}^N \ln \left( \sum_{i=1}^S \exp(\text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i)) \right),$$

where  $\text{sim}(\mathbf{x}_j, \boldsymbol{\mu}_i) = (\mathbf{W}_K \mathbf{x}_j)^\top (\mathbf{W}_Q \boldsymbol{\mu}_i)$ .

$$\mathbf{A}_{j,i} = (\mathbf{W}_K \mathbf{x}_j)^\top (\mathbf{W}_Q \boldsymbol{\mu}_i) \in \mathbb{R}^{N \times S}.$$

**Gradient.** Taking the derivative w.r.t. a single slot  $\boldsymbol{\mu}_i$  yields

$$-\frac{\partial E^{\text{Slot}}}{\partial \boldsymbol{\mu}_i} = \sum_{j=1}^N \text{softmax}_i(\mathbf{A}_{j,i}) \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j.$$

Hence each slot  $\boldsymbol{\mu}_i$  aggregates information from all tokens  $j$ , but the weight is proportional to  $\exp(\mathbf{A}_{j,i})$  normalized across *the slots*  $i$ . This produces the usual iterative update rule:

$$\boldsymbol{\mu}_i^* = \sum_{j=1}^N \text{softmax}_i \left( \boldsymbol{\mu}_i^\top \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j \right) \mathbf{W}_Q^\top \mathbf{W}_K \mathbf{x}_j,$$

which is sometimes referred to as *inverted cross attention*.

## A Proof of the Gradient Decomposition

For completeness, we provide a short derivation of the gradient expression. Recall that our energy is

$$E(\{\mathbf{v}_a\}) = - \sum_{c=1}^A \ln \left( \sum_{\neg \in \mathcal{P}(c)} \exp(\text{sim}(\mathbf{v}_c, \mathbf{v}_{\neg})) \right).$$

Differentiating w.r.t.  $\mathbf{v}_a$ :

$$\frac{\partial E}{\partial \mathbf{v}_a} = - \sum_{c=1}^A \frac{\partial}{\partial \mathbf{v}_a} \ln \left( \sum_{\neg \in \mathcal{P}(c)} \exp(\text{sim}(\mathbf{v}_c, \mathbf{v}_{\neg})) \right).$$

Inside the sum, only terms  $\text{sim}(\mathbf{v}_c, \mathbf{v}_{\neg})$  with  $\neg = a$  or  $c = a$  will contribute. Carefully extracting these leads to the “being explained by parents” plus “explaining children” split in the main text:

$$-\frac{\partial E}{\partial \mathbf{v}_a} = \sum_{\neg \in \mathcal{P}(a)} \text{softmax}_{\neg}(\text{sim}(\mathbf{v}_a, \mathbf{v}_{\neg})) \frac{\partial \text{sim}(\mathbf{v}_a, \mathbf{v}_{\neg})}{\partial \mathbf{v}_a} + \sum_{\substack{c=1 \\ a \in \mathcal{P}(c)}}^A \text{softmax}_a(\text{sim}(\mathbf{v}_c, \mathbf{v}_a)) \frac{\partial \text{sim}(\mathbf{v}_c, \mathbf{v}_a)}{\partial \mathbf{v}_a}.$$