

1 Linear Algebra

Column wise decomposition. Any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ can be decomposed into the sum of its columns:

$$\mathbf{A} = \sum_{j=1}^n \mathbf{A}_{:j} e_j^\top, \quad (1)$$

where e_j are standard basis vectors of \mathbb{R}^n . Notice that this is a rank 1 decomposition.

Row wise decomposition. Any matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ can be decomposed into the sum of its rows:

$$\mathbf{A} = \sum_{i=1}^m e_i \mathbf{A}_{i:}^\top, \quad (2)$$

where e_i are standard basis vectors of \mathbb{R}^m . Notice that this is a rank 1 decomposition.

2 LLM Training

Scaling the logits after LLM head. We usually apply RMS norm to normalize (along the last dimension E , where it stands for model dimension, B means batch size and L means sequence length) the tensor $\mathbf{X} \in \mathbb{R}^{B \times L \times E}$ we feed into LLM head, and obtain the corresponding logits l . After RMS normalization, each tensor corresponding to the token $x_i \in \mathbb{R}^E$ will then have $\text{RMS}(x_t) = 1$. Now notice that for each coordinate $x_{i,t}$, $t \in [E]$, treating as a random variable, its variance is given by

$$\text{Var}(x_{i,t}) = \mathbb{E}[x_{i,t}^2] - (\mathbb{E}[x_{i,t}])^2, \quad (3)$$

and if it is zero-mean (or small), then $\text{Var}(x_{i,t}) \simeq \mathbb{E}[x_{i,t}^2]$, which is to say that second moment reflects the variance.

The next step is to use the empirical observation that for linear layers, hidden vectors tend to be approximatedly rotation-invariant (isotropic), i.e., each coordinate behaves like the others, so we can use the second moment over the coordinate in a token to replace the actual second moment. And the former, is given by

$$\text{Var}(x_{i,t}) \simeq \frac{1}{E} \sum_{t=1}^E x_{i,t}^2 = 1. \quad (4)$$

Now we start to consider the logits, which is generated by

$$l_{j,i} = w_j^\top x_i = \sum_{t=1}^E w_{j,t} x_{i,t}.$$

If we assume each weight entry $w_{j,t}$ are i.i.d. with variance σ^2 the logits variance is give by

$$\text{Var}(l_{j,i}) = \sum_{t=1}^E \sigma^2 \text{Var}(x_{t,i}) \simeq E\sigma^2.$$

So the standard deviation $\sim \sqrt{E}$. To ensure that logits do not scale with the model dimension, we scale it by \sqrt{E} .