

A APPENDIX

Here we provide the derivation for Equation 1. Consider the Taylor expansion of this expression around $\rho = 0$

$$\|\text{Swish}(\mathbf{x})\|^2 = \frac{1}{4}\|\mathbf{x}\|^2 + O(|\rho|). \quad (8)$$

Putting this together with the mean squared norm of the feature embedding

$$E[\|\text{Emb}(\mathbf{x})\|^2] = E[\|\text{Swish}(\gamma \text{BatchNorm}(\text{PreEmb}(\mathbf{x})) + \beta)\|^2], \quad (9)$$

we have

$$E[\|\text{Emb}(\mathbf{x})\|^2] = \frac{1}{4}E[\|\gamma \text{BatchNorm}(\text{PreEmb}(\mathbf{x})) + \beta\|^2] + O(|\rho|). \quad (10)$$

This further simplifies due to the normalizing effect of BatchNorm

$$E[\|\gamma \text{BatchNorm}(\text{PreEmb}(\mathbf{x})) + \beta\|^2] = \sum_i E[(\gamma_i \text{BatchNorm}(\text{PreEmb}(\mathbf{x}))_i + \beta_i)^2] \quad (11)$$

$$= \|\gamma\|^2 + \|\beta\|^2, \quad (12)$$

yielding the result

$$E[\|\text{Emb}(\mathbf{x})\|^2] = \frac{1}{4}\|\gamma\|^2 + \frac{1}{4}\|\beta\|^2 + O(|\rho|). \quad (13)$$