

Theorem 1. *For the canonical MambaZero model with dimensions $d = N = 2^{k+1} = 4$, $e = 1$, and convolution window $w = 2$, there is a choice of parameters such that the model prediction is exactly equal to the Laplacian estimator for random first-order Markov chains. More formally, for any $\beta > 0$, there exists a set of parameters θ such that, for all sequences $(x_t)_{t \geq 1}$ and all $t \geq 1$,*

$$D_{\text{KL}} \left(\mathbb{P}_\beta^{(1)}(\cdot \mid x_1^t) \parallel \mathbb{P}_\theta(\cdot \mid x_1^t) \right) = 0.$$

Proof. Let $\beta > 0$ be the constant of the considered add-constant estimator. Let us fix $a = 0$ and $\Delta_t = 1$, so that $a_t = 1$, for all $t \geq 1$. This can be done by picking, e.g., $\mathbf{w}_\Delta = \mathbf{0}$ and δ such that $\text{softplus}(\delta) = 1$. Note that one can

$$\text{conv}_X = \begin{pmatrix} \alpha_{00} & \alpha_{01} \\ \alpha_{10} & \alpha_{11} \end{pmatrix}, \quad \text{conv}_B = \begin{pmatrix} \gamma_{00} & \gamma_{01} \\ \gamma_{10} & \gamma_{11} \end{pmatrix} \quad (1)$$

where each row corresponds to the kernel weights applied time-wise to each coordinate of the input sequence $(\mathbf{x}_t)_{t \geq 1}$. Since the window for conv_C is $w_C = 1$, we can simply assume w.l.o.g. that $C_t = W_C \mathbf{x}_t$.

Let us take the embedding vectors to be $\mathbf{e}_0 = (1, 0, 0, 0)^\top$ and $\mathbf{e}_1 = (0, 1, 0, 0)^\top$, and take

$$W_X = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}. \quad (2)$$

Take also

$$\text{conv}_X = \begin{pmatrix} 1 & 1 \\ 3 & -1 \\ 1 & 1 \\ 3 & -1 \end{pmatrix}, \quad (3)$$

so that one has, after the application of W_X and conv_X ,

$$X^{(00)} = \begin{pmatrix} 2 \\ 2 \\ 0 \\ 0 \end{pmatrix}, \quad X^{(01)} = \begin{pmatrix} 1 \\ 3 \\ 1 \\ -1 \end{pmatrix}, \quad X^{(10)} = \begin{pmatrix} 1 \\ -1 \\ 1 \\ 3 \end{pmatrix}, \quad X^{(11)} = \begin{pmatrix} 0 \\ 0 \\ 2 \\ 2 \end{pmatrix}. \quad (4)$$

Take also $\text{conv}_B = \text{conv}_X$ and $W_B = W_X$, so that

$$B_0 = W_B \mathbf{e}_0 = \begin{pmatrix} 1 \\ 1 \\ 0 \\ 0 \end{pmatrix}, \quad B_1 = W_B \mathbf{e}_1 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}, \quad (5)$$

and take $W_C = \frac{1}{4}W_B$, so that $C^{(0)} = \frac{1}{4}B_0$ and $C^{(1)} = \frac{1}{4}B_1$.

The final logit vector is in general equal to

$$\text{logit}_t = W_\ell \mathbf{x}_t + W_\ell W_o X^{(x_1)} B^{(x_1)^\top} C^{(x_t)} + \sum_{ij} n_{ij} W_\ell W_o X^{(ij)} B^{(ij)^\top} C^{(x_t)}. \quad (6)$$

Note that the application of convolution to a given sequence of vectors \mathbf{z}_1^t can be rewritten as a linear matrix-form operation. For example, for conv_X , one has that

$$\text{conv}_X(\mathbf{z}_t) = D_X^{(0)} \mathbf{z}_{t-1} + D_X^{(1)} \mathbf{z}_t \quad (7)$$

where $D_X^{(0)} = \text{diag}(\alpha_{00}, \alpha_{10})$ and $D_X^{(1)} = \text{diag}(\alpha_{01}, \alpha_{11})$ are diagonal matrices. The same holds for conv_B , with corresponding diagonal matrices $D_B^{(0)}$ and $D_B^{(1)}$. Using this fact, we can rewrite the logit formula as

$$\begin{aligned} \text{logit}_t = & W_\ell \mathbf{x}_t + W_\ell W_o D_X^{(1)} X_{x_1} B_{x_1}^\top D_B^{(1)} C^{(x_t)} \\ & + \sum_{ij} n_{ij} W_\ell W_o (D_X^{(0)} X_i + D_X^{(1)} X_j) (B_i^\top D_B^{(0)} C^{(x_t)} + B_j^\top D_B^{(1)} C^{(x_t)}). \end{aligned} \quad (8)$$

Note that, with the choice of parameters above, one has

$$D_X^{(1)} X_{x_1} B_{x_1}^\top D_B^{(1)} C^{(i)} = \mathbf{0} \quad (9)$$

and also,

$$B_0^\top D_B^{(0)} C^{(1)} + B_0^\top D_B^{(1)} C^{(1)} = 0 \quad (10)$$

$$B_0^\top D_B^{(0)} C^{(1)} + B_1^\top D_B^{(1)} C^{(1)} = 0 \quad (11)$$

$$B_1^\top D_B^{(0)} C^{(0)} + B_0^\top D_B^{(1)} C^{(0)} = 0 \quad (12)$$

$$B_1^\top D_B^{(0)} C^{(0)} + B_1^\top D_B^{(1)} C^{(0)} = 0, \quad (13)$$

that is, only the desired counts are kept in the final logit, depending on the current symbol. Furthermore, for the relevant counts, one has

$$B_0^\top D_B^{(0)} C^{(0)} + B_0^\top D_B^{(1)} C^{(0)} = 1 \quad (14)$$

$$B_0^\top D_B^{(0)} C^{(0)} + B_1^\top D_B^{(1)} C^{(0)} = 1 \quad (15)$$

$$B_1^\top D_B^{(0)} C^{(1)} + B_0^\top D_B^{(1)} C^{(1)} = 1 \quad (16)$$

$$B_1^\top D_B^{(0)} C^{(1)} + B_1^\top D_B^{(1)} C^{(1)} = 1. \quad (17)$$

Hence, the final logit becomes

$$\text{logit}_t = W_\ell \mathbf{e}_0 + \sum_j n_{0j} W_\ell W_o (D_X^{(0)} X_0 + D_X^{(1)} X_j) \quad (18)$$

for $x_t = 0$, and

$$\text{logit}_t = W_\ell \mathbf{e}_1 + \sum_j n_{1j} W_\ell W_o (D_X^{(0)} X_1 + D_X^{(1)} X_j) \quad (19)$$

for $x_t = 1$. Finally, take

$$W_\ell = \begin{pmatrix} \beta & \beta & 1 & 0 \\ \beta & \beta & 0 & 1 \end{pmatrix} \quad (20)$$

and

$$W_o = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 1 & -\frac{1}{2} & \frac{1}{4} & -\frac{1}{4} \\ \frac{1}{4} & -\frac{1}{4} & 1 & -\frac{1}{2} \end{pmatrix}. \quad (21)$$

With this choice, we get, for all $t \geq 1$,

$$\text{logit}_t = \begin{pmatrix} \beta \\ \beta \end{pmatrix} + n_{x_t,0} \begin{pmatrix} 1 \\ 0 \end{pmatrix} + n_{x_t,1} \begin{pmatrix} 0 \\ 1 \end{pmatrix}. \quad (22)$$

After the normalization, we finally get

$$f_{\boldsymbol{\theta}}(x_1^t) = \left(\frac{n_{00} + \beta}{n_{00} + n_{01} + 2\beta}, \frac{n_{01} + \beta}{n_{00} + n_{01} + 2\beta} \right)^\top \quad (23)$$

if $x_t = 0$, and

$$f_{\boldsymbol{\theta}}(x_1^t) = \left(\frac{n_{10} + \beta}{n_{10} + n_{11} + 2\beta}, \frac{n_{11} + \beta}{n_{10} + n_{11} + 2\beta} \right)^\top \quad (24)$$

if $x_t = 1$. This is precisely the required add- β Laplacian estimator. \square