

weight-normalization VS threshold-balancing¹

2022 年 4 月

符号表示

l : $l \in \{1, \dots, L\}$ a network with L layers, layer l
 W_{ij}^l : weight connection between neuron i in layer l and neuron j in layer $l-1$
 b_i^l : neuron i bias in layer l
 a_i^l : neuron i activation value(after relu) in layer l
 M^l : the number of neurons in layer l
 $V_i^l(t)$: neuron i membrane potential in layer l and time t
 $z_i^l(t)$: neuron i inputs in layer l and time t
 $\Theta_{t,i}^l$: a step function indicating the occurrence of a spike at time t
 $N_i^l(t)$: the number of spike generated in neuron i layer l for a total time t
 $r_i^l(t)$: firing rate of neuron i in layer l for a total time t
 V_{thr}^l : voltage threshold in layer l

满足关系

$$\Theta_{t,i}^l = \Theta \left(V_i^l(t-1) + z_i^l(t) - V_{\text{thr}} \right), \text{ with } \Theta(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{else} \end{cases}$$

$$N_i^l(t) = \sum_{t'=1}^t \Theta_{t',i}^l$$

$$r_i^l(t) = N_i^l(t)/t$$

$$a_i^l = \sum_{j=1}^{M^{l-1}} W_{ij}^l a_j^{l-1} + b_i^l$$

$$V_i^l(t) = V_i^l(t-1) + z_i^l(t) - V_{\text{thr}} \Theta_{t,i}^l$$

$$z_i^l(t) = V_{\text{thr}} \left(\sum_j^{M^{l-1}} W_{ij}^l \Theta_{t,i}^l + b_i^l \right)$$

*

¹*: 下面探讨针对的都是软重置的脉冲发放方式

对第一层

a) 权重归一化:

$$V_i^l(t) = V_i^l(t-1) + \underbrace{\sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l}_{a_i^l} - \Theta_{t,i}^l$$

移项并对 T 个时间步长求和:

$$\sum_{t=0}^T \Theta_{t,i}^l = \sum_{t=0}^T \sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l - V_i^l(T)$$

由于输入的 x 是固定值, 所以第一层的 $Wx+b$ 是常数项, 式子变成

$$\sum_{t=0}^T \Theta_{t,i}^l = \left(\sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l \right) T - V_i^l(T)$$

两边除以 T

$$r_i^l(T) = \underbrace{\sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l}_{a_i^l} - \frac{V_i^l(T)}{T} \quad (1)$$

Diehl 当时对权重归一化的解释为: 转换后由于阈值和权重的比例不合适, 导致神经元欠激活或者过激活, 提出权重归一化来减少这个误差。实际上由式子 (1), **另一种解释方法为:** 要想 RELU 单元用 firing rate approximation of an IF neuron with no refractory period 近似, 这里的 firing rate 即 $r_i^l(T) \in [0, 1]$, 需要将模拟激活值 a_i^l 映射到 $[0, 1]$, 才能用不带不应期的 IF 神经元脉冲发放率做近似。这就是为什么除以激活值最大值, 即权重归一化的数学解释

同时只有对正的输入, 膜电势才会积累到超过阈值发放脉冲, 这就使得负的输出膜电势越来越负, 脉冲发放率为零, 完美地用 IF 神经元等价了 RELU 激活函数。

所以注意 (1) 式子中的 W, b 是归一化后的值, 若是换一种表达, 即如果 W, b 是 ANN 的原始权重, 那么 (1) 式子应该写为:

$$r_i^l(T) = \underbrace{\sum_j^{M^{l-1}} \frac{W_{ij}^l}{a_{max}^l} x_j + \frac{b_i^l}{a_{max}^l}}_{a_i^l/a_{max}^l} - \frac{V_i^l(T)}{T} \quad (2)$$

b) 阈值平衡:

上周讨论了发放脉冲是 1 还是 V_{thr} 的问题, 显然发放 V_{thr} 从生物上不符合脉冲的特性, 但是下面我们先说明: 发放 V_{thr} 和权重归一化等价。这时把 $V_{thr} \Theta_{t,i}^l$ 看做一个脉冲数

$$V_i^l(t) = V_i^l(t-1) + \underbrace{\sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l}_{a_i^l} - V_{thr} \Theta_{t,i}^l \quad (3)$$

移项并对 T 个时间步长求和：

$$V_{\text{thr}} \sum_{t=0}^T \Theta_{t,i}^l = \sum_{t=0}^T \sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l - V_i^l(T)$$

由于输入的 x 是固定值, 所以第一层的 $Wx+b$ 是常数项, 式子变成

$$V_{\text{thr}} \sum_{t=0}^T \Theta_{t,i}^l = \left(\sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l \right) T - V_i^l(T)$$

两边除以 T , 由于 $r_i^l(T) = \frac{N_i^l(T)}{T} = \frac{V_{\text{thr}} \Theta_{t,i}^l}{T}$

$$r_i^l(T) = \underbrace{\sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l}_{a_i^l} - \frac{V_i^l(T)}{T} \quad (4)$$

这里的 $V_i^l(T)$ 是 (2) 式中的 V_{thr} 倍, 在 V_{thr} 设为 a_{max}^l 的情况下, 显然权重归一化后的脉冲发放率乘以最大值后, 等于阈值平衡的脉冲发放率。

————以下先忽略————

注意满足关系一节中的 (*) 式, 它的输入很奇怪, 是输入又乘以一个电压阈值。在发放脉冲为 1 的情况下, (*) 式这样的输入才能使两种不同发放的阈值平衡等价。所以所有输入前面多乘了一项 V_{thr} , (3) 式变换为

$$V_i^l(t) = V_i^l(t-1) + V_{\text{thr}} \underbrace{\sum_j^{M^{l-1}} W_{ij}^l x_j + b_i^l}_{a_i^l} - V_{\text{thr}} \Theta_{t,i}^l \quad (5)$$

由于 $r_i^l(T) = \frac{N_i^l(T)}{T} = \frac{\Theta_{t,i}^l}{T}$, 可以自然得到与 (4) 式一样的结果

对更高层

$$V_i^l(t) = V_i^l(t-1) + \sum_j^{M^{l-1}} W_{ij}^l \Theta_{t,j}^{l-1} + b_i^l - V_{\text{thr}} \Theta_{t,i}^l$$

同样地, 进行 T 次求和:

$$V_{\text{thr}} \sum_{t=0}^T \Theta_{t,i}^l = \sum_{t=0}^T \left(\sum_j^{M^{l-1}} W_{ij}^l \Theta_{t,j}^{l-1} + b_i^l \right) - V_i^l(T)$$

两边同时除以 T

$$r_i^l(T) = \sum_j^{M^{l-1}} W_{ij}^l r_j^{l-1} + b_i^l - \frac{V_i^l(T)}{T} \quad (6)$$

这是阈值平衡的，显然权重归一化后的脉冲发放率乘以最大值后，仍然等于阈值平衡的脉冲发放率，即两者在高层也完全等价。

激活值与脉冲发放率之间的关系

由 (6) 式可得不同层间的脉冲发放率关系。为了方便表示，我们用权重归一化的形式进行展开。

$$\begin{aligned}
r_i^l(T) &= \sum_j^{M^{l-1}} \frac{W_{ij}^l a_{max}^{l-1}}{a_{max}^l} r_j^{l-1} + \frac{b_i^l}{a_{max}^l} - \frac{V_i^l(T)}{T} \\
&= \sum_j^{l-1} \sum_j^{M^{l-1}} \frac{W_{ij}^l a_{max}^{l-1}}{a_{max}^l} \left(\sum_k^{M^{l-2}} \frac{W_{jk}^{l-1} a_{max}^{l-2}}{a_{max}^{l-1}} r_k^{l-2} + \frac{b_j^{l-1}}{a_{max}^{l-1}} - \frac{V_j^{l-1}(T)}{T} \right) + \frac{b_i^l}{a_{max}^l} - \frac{V_i^l(T)}{T} \\
&= \sum_j^{l-1} \sum_j^{M^{l-1}} \frac{W_{ij}^l a_{max}^{l-1}}{a_{max}^l} \left(\sum_k^{M^{l-2}} \frac{W_{jk}^{l-1} a_{max}^{l-2}}{a_{max}^{l-1}} \left(\sum \dots \sum_m^{M^1} \underbrace{\frac{W_{1m}^1}{a_{max}^1} x_m + \frac{b_m^1}{a_{max}^1}}_{a_m^1/a_{max}^1} - \frac{V_i^1(T)}{T} + \dots + \frac{b}{a_{max}} - \frac{V_i^l(T)}{T} \right) \right. \\
&\quad \left. + \frac{b_j^{l-1}}{a_{max}^{l-1}} - \frac{V_j^{l-1}(T)}{T} \right) + \frac{b_i^l}{a_{max}^l} - \frac{V_i^l(T)}{T}
\end{aligned}$$

用 ΔV_i^l 表示 $\frac{V_i^l(T)}{T}$

则上面的式子变为

$$r_i^l(T) = \frac{a_i^l}{a_{max}^l} - \Delta V_i^l - \sum_j^{M^{l-1}} \frac{W_{ij}^l a_{max}^{l-1}}{a_{max}^l} \Delta V_j^l - \dots - \sum_j^{M^{l-1}} \frac{W_{ij}^l a_{max}^{l-1}}{a_{max}^l} \dots \sum_k^{M^1} \frac{W_{1k}^2 a_{max}^1}{a_{max}^2} \Delta V_k^1 \quad (7)$$

所以从理论上讲，用脉冲发放率 * 最后一层的最大激活值来近似模拟激活值 a_i^l ，有 (7) 式中，减号后面一长串的固有误差，这个误差每项都乘了 ΔV ，即 $\frac{V(T)}{T}$ 。所以增加时间步长可以有效地减少转换的损失。对于分类问题，虽然这个误差对每个神经元都存在，但是时间步长的增大，这个误差往往不大，所以对取最大值作为分类结果这样的任务，需要的时间步长不大。但对于检测，这样的误差会引起浮点数计算的改变，所以需要更大的时间步长来减少这一固有损失。

那么既然有了固有误差，人为地把它补上即把在得到脉冲发放率后，将右边误差加到发放率上，就可以无损地还原模拟激活值了，而且不受时间步长的影响。遗憾的是理论上确实可行，但是这完全没有生物合理性，也失去了脉冲神经网络的意义；从代码实现来讲，这些权重是使得输入 $W\Theta + b$ 为正的权重，将每层的这些权重找出来再和 ΔV 做乘法就过于复杂了。所以在为了简化，转换中人们还是带着这样的误差在做视觉任务。