# RNN Diagram - ReLU Activation on Hidden Layer
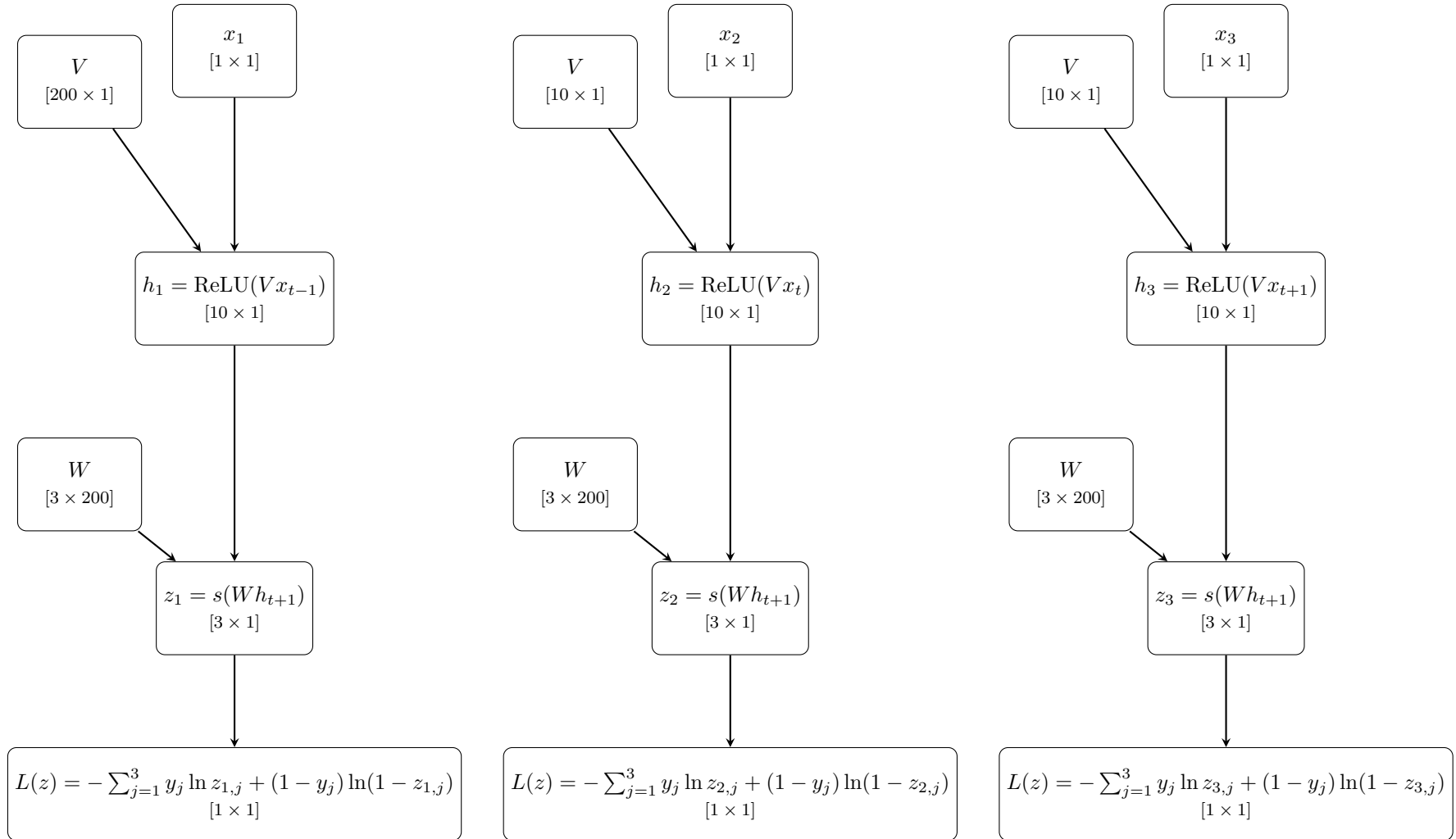
The constructed neural net follows the diagram below. $x \in \mathbb{R}^{24}$, a feature vector with (4 hours $\times$ 6 intervals) = 24 features corresponding to occupancies every 10 minutes for 4 hours before the current time. For clarity, shapes of the solutions are given in brackets. Subscript $j$ denotes a row index. $A_j^T$ indicates the transpose of $A_j$.

| | | |
|---|---|---|
| $V$ $[200 \times 1]$ $\quad$ $x_1$ $[1 \times 1]$ | $V$ $[10 \times 1]$ $\quad$ $x_2$ $[1 \times 1]$ | $V$ $[10 \times 1]$ $\quad$ $x_3$ $[1 \times 1]$ |
| $h_1 = \text{ReLU}(Vx_{t-1})$ $[10 \times 1]$ | $h_2 = \text{ReLU}(Vx_t)$ $[10 \times 1]$ | $h_3 = \text{ReLU}(Vx_{t+1})$ $[10 \times 1]$ |
| $W$ $[3 \times 200]$ | $W$ $[3 \times 200]$ | $W$ $[3 \times 200]$ |
| $z_1 = s(Wh_{t+1})$ $[3 \times 1]$ | $z_2 = s(Wh_{t+1})$ $[3 \times 1]$ | $z_3 = s(Wh_{t+1})$ $[3 \times 1]$ |
| $L(z) = -\sum_{j=1}^{3} y_j \ln z_{1,j} + (1 - y_j)\ln(1 - z_{1,j})$ $[1 \times 1]$ | $L(z) = -\sum_{j=1}^{3} y_j \ln z_{2,j} + (1 - y_j)\ln(1 - z_{2,j})$ $[1 \times 1]$ | $L(z) = -\sum_{j=1}^{3} y_j \ln z_{3,j} + (1 - y_j)\ln(1 - z_{3,j})$ $[1 \times 1]$ |

$A_k = V x_k$

$\nabla_V L =$
$(\nabla_{z_f} L)(\nabla_V z_{t+1}) =$
$(\nabla_{z_f} L)(\nabla_{h_{t+1}} z_{t+1})(\nabla_V h_{t+1}) =$
$\left(\frac{1-y_j}{1-z_j} - \frac{y_j}{z_j}\right)\left(z_j(1-z_j)W_j^T\right)(\nabla_V h_{t+1}) =$
$(z_j - y_j)W_j^T(\nabla_V h_{t+1}) =$
$W_j^T(z_j - y_j)(\nabla_V h_{t+1}) =$
$W^T(z-y)(\nabla_V h_{t+1}) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})(\nabla_V A_{t+1}) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})(\nabla_V V x_{t+1}) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})(\nabla_V U h_t) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})(\nabla_V V x_{t+1}) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})(U\nabla_V h_t) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})\left(I^{[200\times200]} x_{t+1}\right) =$

$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t)(\nabla_V A_t) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})\left(I^{[200\times200]} x_{t+1}\right) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t)(\nabla_V U h_{t-1} + \nabla_V V x_t) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})\left(I^{[200\times200]} x_{t+1}\right) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t)(\nabla_V U h_{t-1}) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t)(\nabla_V V x_t) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})\left(I^{[200\times200]} x_{t+1}\right) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t) U (\nabla_V h_{t-1}) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U \left(I^{[200\times200]} x_t\right)(\nabla_{A_t} h_t) + W^T(z-y)\left(I^{[200\times200]} x_{t+1}\right)(\nabla_{A_{t+1}} h_{t+1}) =$

$\nabla_U L =$
$(\nabla_{z_{t+1}} L)(\nabla_U z_{t+1}) =$
$(\nabla_{z_{t+1}} L)(\nabla_{h_{t+1}} z_{t+1})(\nabla_U h_{t+1}) =$
$\left(\frac{1-y_j}{1-z_j} - \frac{y_j}{z_j}\right)\left(z_j(1-z_j)W_j^T\right)(\nabla_U h_{t+1}) =$
$(z_j - y_j)W_j^T(\nabla_U h_{t+1}) =$
$W_j^T(z_j - y_j)(\nabla_U h_{t+1}) =$
$W^T(z-y)(\nabla_U h_{t+1}) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})(\nabla_U A_{t+1}) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})(\nabla_U U h_t) =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1})[U\nabla_U h_t + h_t] =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_U h_t) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) h_t =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t)(\nabla_U A_t) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) h_t =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t)(\nabla_U U h_{t-1}) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) h_t =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t)[U\nabla_U h_{t-1} + h_{t-1}] + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) h_t =$
$W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t) U (\nabla_U h_{t-1}) + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) U (\nabla_{A_t} h_t) h_{t-1} + W^T(z-y)(\nabla_{A_{t+1}} h_{t+1}) h_t =$

$\nabla_W L =$
$\left(\nabla_{z_f} L\right)\left(\nabla_W z_f\right)$
$\left(\nabla_{z_j} L\right)\left(\nabla_{W_j} z_j\right)$
$\left(\frac{1-y_j}{1-z_j} - \frac{y_j}{z_j}\right) z_j(1-z_j)h_f^T$
$(z_j - y_j)h_f^T$
$(z - y)h_f^T$ $\quad$ [3 × 200]

$\nabla_{z_j} L = \frac{1-y_j}{1-z_j} - \frac{y_j}{z_j}$
$\nabla_{W_j} z_j = z_j(1-z_j)h^T$ $\quad$ [1 × 200]
$\nabla_h z_j = z_j(1-z_j)W_j^T$ $\quad$ [201 × 1]
$\nabla_A h = \text{(elementwise)} \begin{cases} 1, & h_j > 0 \\ 0, & h_j < 0 \end{cases}$
$\nabla_U h = \text{recursive expansion}$
$\nabla_U h = \text{recursive expansion}$

We can substitute the following expressions:
$\quad \nabla_h z_j = z_j(1-z_j)W_j^T$ $\quad$ [201 × 1]
$\quad \nabla_{V_j} h_j = \text{sech}^2(V_j x)x^T$ $\quad$ [1 × 785]

We can also use backpropagation to show:
$\quad \nabla_{W_j} L = \left(\frac{1-y_j}{1-z_j} - \frac{y_j}{z_j}\right)\left(z_j(1-z_j)h^T\right)$ $\quad$ [1 × 201]
$\quad \nabla_h L = \sum_{j=1}^{26} \left(\frac{1-y_j}{1-z_j} - \frac{y_j}{z_j}\right)\left(z_j(1-z_j)W_j^T\right)$ $\quad$ [201 × 1]
$\quad \nabla_{V_j} L = \left(\nabla_h L\right)_j \text{sech}^2(V_j x)x^T$ $\quad$ [1 × 785]

To enhance calculation efficiency, we can reduce some of these equations into matrices. First, let $\mathcal{Q}$ be the matrix

$$\mathcal{Q} = \begin{bmatrix} z_1(1-z_1)\left(\frac{1-y_1}{1-z_1} - \frac{y_1}{z_1}\right) \\ z_2(1-z_2)\left(\frac{1-y_2}{1-z_2} - \frac{y_2}{z_2}\right) \\ \vdots \\ z_{26}(1-z_{26})\left(\frac{1-y_{26}}{1-z_{26}} - \frac{y_{26}}{z_{26}}\right) \end{bmatrix} = \begin{bmatrix} z_1 - y_1 \\ z_2 - y_2 \\ \vdots \\ z_{26} - y_{26} \end{bmatrix}$$

With this, we can reexpress the above equations.

(1)

$$\nabla_{W_j} L = \mathcal{Q}_j h^T \qquad [1 \times 201]$$

or

$$\nabla_W L = \begin{bmatrix} \mathcal{Q}_1 h^T \\ \mathcal{Q}_2 h^T \\ \vdots \\ \mathcal{Q}_{26} h^T \end{bmatrix} = \mathcal{Q} h^T = \mathcal{Q} \otimes h \qquad [26 \times 201]$$

(2)

$$\nabla_h L = \sum_{j=1}^{26} \mathcal{Q}_j W_j^T \qquad [201 \times 1]$$

or

$$\nabla_h L = \left[ \mathcal{Q}_1 W_1^T + \mathcal{Q}_2 W_2^T + ... + \mathcal{Q}_{26} W_{26}^T \right]$$

$$\nabla_h L = \left[ W_1^T \mathcal{Q}_1 + W_2^T \mathcal{Q}_2 + ... + W_{26}^T \mathcal{Q}_{26} \right]$$

$$\nabla_h L = \begin{bmatrix} W_1^T & W_2^T & \dots & W_{26}^T \end{bmatrix} \begin{bmatrix} \mathcal{Q}_1 \\ \mathcal{Q}_2 \\ \vdots \\ \mathcal{Q}_{26} \end{bmatrix} = W^T \mathcal{Q} \qquad [201 \times 1]$$

(3)

Additionally, let

$$\mathcal{S} = \begin{bmatrix} \operatorname{sech}^2(V_1 x) \\ \operatorname{sech}^2(V_2 x) \\ \vdots \\ \operatorname{sech}^2(V_{200} x) \end{bmatrix} = \operatorname{sech}^2(Vx)$$

Then, $\nabla_{V_j} L = (W^T \mathcal{Q})_j \operatorname{sech}^2(V_j x) x^T \qquad [1 \times 785]$

or

$$\nabla_V L = \begin{bmatrix} (W^T \mathcal{Q})_1 \operatorname{sech}^2(V_1 x) x^T \\ (W^T \mathcal{Q})_2 \operatorname{sech}^2(V_2 x) x^T \\ \vdots \\ (W^T \mathcal{Q})_{200} \operatorname{sech}^2(V_{200} x) x^T \end{bmatrix} = \begin{bmatrix} (W^T \mathcal{Q})_1 \operatorname{sech}^2(V_1 x) \\ (W^T \mathcal{Q})_2 \operatorname{sech}^2(V_2 x) \\ \vdots \\ (W^T \mathcal{Q})_{200} \operatorname{sech}^2(V_{200} x) \end{bmatrix} x^T$$

$$\nabla_V L = \left( \begin{bmatrix} (W^T \mathcal{Q})_1 \\ (W^T \mathcal{Q})_2 \\ \vdots \\ (W^T \mathcal{Q})_{200} \end{bmatrix} \circ \begin{bmatrix} \operatorname{sech}^2(V_1 x) \\ \operatorname{sech}^2(V_2 x) \\ \vdots \\ \operatorname{sech}^2(V_{200} x) \end{bmatrix} \right) x^T = (W^T \mathcal{Q}) \circ \mathcal{S} x^T$$

Since we are updating our matrices $V$ and $W$ using stochastic gradient descent, we repeat the following process:

$V, W \leftarrow$ weight matrices initialized randomly from normal distribution with mean $\mu = 0$ and $\sigma^2 = (...)$
while (continue = True or $L(z) > 0$)
      Forward calculation $[\ h = \tanh(Vx)\ \ \rightarrow\ \ z = s(Wh)\ \ \rightarrow\ \ L(z)\ ]$
      Backward calculation to return $\nabla_V L$ and $\nabla_W L$
      $V \leftarrow V - \epsilon \nabla_V L)$
      $W \leftarrow W - \epsilon \nabla_W L$
return $V, W$

where, using our derived equations from above, the update rules are more specifically
$$V \leftarrow V - \epsilon (W^T \mathcal{Q}) \circ \mathcal{S} x^T$$
$$W \leftarrow W - \epsilon (\mathcal{Q} \otimes h).$$

$h = \text{ReLU}(U h_{t-1} + V x_t)$
$\partial h / \partial h_{t-1} = \text{RLD}(U h_{t-1} + V x_t) U^T$
$\partial h / \partial V = \text{RLD}(U h_{t-1} + V x_t)(\partial(U h_{t-1})/\partial V + \partial(V x_t)/\partial V) = \text{RLD}(U h_{t-1} + V x_t)(U^T \frac{\partial h_{t-1}}{\partial h_{t-2}}(...)\frac{\partial h_0}{\partial V} + I x_t)$
$\partial h_0 / \partial V =$
$\partial h / \partial U = \text{RLD}(U h_{t-1} + V x_t)(\partial(U h_{t-1})/\partial U + \partial(V x_t)/\partial U) = \text{RLD}(U h_{t-1} + V x_t)(U^T \frac{\partial h_{t-1}}{\partial h_{t-2}}(...)\frac{\partial h_0}{\partial V} + I x_t)$

$y = \sin(Ax)$
$\partial y / \partial x = \frac{\partial(Ax)}{\partial x}\cos(Ax) = A^T \cos(Ax)$
$[200 \times 200] = [200 \times 1](...)$