



Feed Forward

- Dot product between input x and input weight U
- Dot product between initial hidden state h_o & weight matrix W
- Summation between output of both dot product operations & bias b , resulting in logits a_t
- Pass a_t into the activation node g which uses the hyperbolic tangent (\tanh) as activation function, resulting in h_t
- Dot product of h_t & V , added with c (this represents a fully connected layer after the new hidden state is calculated), outputting o_t
- Pass o_t to softmax to output \hat{y}_t
- Pass \hat{y}_t to CE_{Loss} node to get L_t , which represents the loss for the branch at time t

$$a_t = Ux_t + Wh_{t-1} + b$$

$$h_t = \tanh(a_t)$$

$$o_t = Vh_t + c$$

$$\hat{y}_t = \text{softmax}(o_t)$$

$$L_t = CE(\hat{y}_t, y_t)$$

Dimensionality

Input

$$X_t: [\text{batch-size}, \text{input-dim}]$$

Weight matrices & Biases

$$U: [\text{input-dim}, \text{hidden-dim}]$$

$$W: [\text{hidden-dim}, \text{hidden-dim}]$$

$$b: [1, \text{hidden-dim}]$$

$$V: [\text{hidden-dim}, \text{output-dim}]$$

$$c: [1, \text{output-dim}]$$

Hidden State

$$h_t: [\text{batch-size}, \text{hidden-dim}]$$

Output Before Softmax

$$o_t: [\text{batch-size}, \text{output-dim}]$$

Predicted output & True output

$$\hat{y}_t \text{ \& } y_t: [\text{batch-size}, \text{output-dim}]$$

Loss

$$L_t: [\text{batch-size}, 1]$$