- Notation and Problem Statement

$$X = (x^1, x^2, \cdots, x^n)^\top \quad n \text{ driving series} : n \text{ features}$$

$$= (X_1, X_2, \cdots, X_T) \in \mathbb{R}^{n \times T}$$

where $x^k = (x_1^k, x_2^k, \cdots, x_T^k)^\top$

$$
\begin{matrix}
x_1^1 & x_1^2 & \cdots & x_1^n \\
x_2^1 & x_2^2 & \cdots & x_2^n \\
\vdots & \vdots & \ddots & \vdots \\
x_T^1 & x_T^2 & \cdots & x_T^n
\end{matrix}
$$

$X_t = (x_t^1, x_t^2, \cdots, x_t^n)^\top \in \mathbb{R}^n$

시간 $t$ 에서 $n$개의

feature 에 대한 행렬

$(y_1, y_2, \cdots, y_{T-1})$ with $y_t \in \mathbb{R}$

$$\hat{y}_T = F(y_1, \cdots, y_{T-1}, X_1, \cdots, X_T)$$

- Encoder with input attention

Input sequence $X = (x_1, x_2, \cdots, x_T)$, $x_t \in \mathbb{R}^n$

(eqn 2)

$lh_t = f_1(lh_{t-1}, x_t)$, $lh_t \in \mathbb{R}^m$   m: size of hidden state

we use an LSTM unit as $f_1$ to capture

$\underline{\hspace{4cm}\text{long-term dependencies}}$

## LSTM unit

$f_t$: forget gate, $i_t$: input gate, $O_t$: output gate

$f_t = \sigma(W_f [lh_{t-1}; x_t] + lb_f)$          eqn 3

$i_t = \sigma(W_i [lh_{t-1}; x_t] + lb_i)$          eqn4

$O_t = \sigma(W_o [lh_{t-1}; x_t] + lb_o)$          eqn 5

$S_t = f_t \odot S_{t-1} + i_t \odot \tanh(W_s [lh_{t-1}; x_t] + lb_s)$

$lh_t = O_t \odot \tanh(S_t)$          eqn 7       $\hookrightarrow$ eqn 6

$[lh_{t-1}; x_t] \in \mathbb{R}^{m+n}$ : concatenation of $lh_{t-1}$ and $x_t$

$W_{f-s} \in \mathbb{R}^{m \times (m+n)}$              $W[lh_{t-1}; x_t] \in \mathbb{R}^m$

$lb_{f-s} \in \mathbb{R}^m$              $\rfloor$ parameters to learn

## Input attention

$$e_t^k = v_e^+ \tanh(W_e [h_{t-1}; s_{t-1}] + U_e x^k) \qquad \text{eqn 8}$$

$$\alpha_t^k = \frac{\exp(e_t^k)}{\sum_{i=1}^n \exp(e_t^i)} \quad : \text{attention weight} \qquad \text{eqn 9}$$

$$\text{(softmax function applied)}$$

where $v_e \in \mathbb{R}^T$, $W_e \in \mathbb{R}^{T \times 2m}$ and $U_e \in \mathbb{R}$

$$x^k \in \mathbb{R}$$

$$\tilde{x}_t = (\alpha_t^1 x_t^1, \alpha_t^2 x_t^2, \cdots, \alpha_t^n x_t^n)^\top \qquad \text{eqn 10}$$

$$h_t = f_1(h_{t-1}, \tilde{x}_t) \qquad \text{eqn 11}$$

$\hookrightarrow$ hidden state at time t

$\square$

- Decoder with temporal attention

<span style="color:blue">Temporal attention</span>

$$l_t^i = v_d^\top \tanh(W_d [dl_{t-1} ; s'_{t-1}] + U_d lh_i) \quad 1 \leq i \leq T \qquad \text{eqn } 12$$

$dl_{t-1} \in \mathbb{R}^p$ : previous decoder hidden state 　(p?)

$s'_{t-1} \in \mathbb{R}^p$ : cell state of the LSTM unit

$[dl_{t-1} ; s'_{t-1}] \in \mathbb{R}^{2p}$

$v_d \in \mathbb{R}^m$, $W_d \in \mathbb{R}^{m \times 2p}$, and $U_d \in \mathbb{R}^{m \times n}$ : parameters to learn

$$\beta_t^i = \frac{\exp(l_t^j)}{\sum_{j=1}^{T} \exp(l_t^j)} \quad ; \text{ attention weight} \qquad \text{eqn } 13$$

$$C_t = \sum_{i=1}^{T} \beta_t^i \, lh_i \quad : \text{context vecotor} \qquad \text{eqn } 14$$

$$\tilde{y}_{t-1} = \tilde{w}^\top [y_{t-1} ; C_{t-1}] + \tilde{b} \qquad \text{eqn } 15$$

where $[y_{t-1} ; C_{t-1}] \in \mathbb{R}^{m+1}$

$$dl_t = f_2(dl_{t-1} , \tilde{y}_{t-1}) \qquad \text{eqn } 16$$

## LSTM unit

$$f'_t = \sigma \left( W'_f \left[ dl_{t-1} ; \tilde{y}_{t-1} \right] + b'_f \right) \qquad \text{eqn 17}$$

$$i'_t = \sigma \left( W'_i \left[ dl_{t-1} ; \tilde{y}_{t-1} \right] + b'_i \right) \qquad \text{eqn 18}$$

$$O'_t = \sigma \left( W'_o \left[ dl_{t-1} ; \tilde{y}_{t-1} \right] + b'_o \right) \qquad \text{eqn 19}$$

$$S'_t = f'_t \odot S'_{t-1} + i'_t \odot \tanh \left( W'_s \left[ dl_{t-1} ; \tilde{y}_{t-1} \right] + b'_s \right)$$

$$dl_t = O'_t \odot \tanh \left( S'_t \right) \qquad \text{eqn 21} \qquad \hookrightarrow \text{eqn 20}$$

$$\left[ dl_{t-1} ; \tilde{y}_{t-1} \right] \in \mathbb{R}^{P+1}$$

$$W'_{f-s} \in \mathbb{R}^{P \times (P+1)} \quad \text{and} \quad b'_{f-s} \in \mathbb{R}^P : \text{parameters to learn}$$

$$\hat{y}_T = F \left( y_1, \cdots ; y_{T-1} , x_1, \cdots , x_T \right)$$

$$\quad = V^t \left( W_y \left[ dl_T ; C_T \right] + b_w \right) + b_v$$

$$O \left( y_T, \hat{y}_T \right) = \frac{1}{N} \sum_{i=1}^{N} \left( \hat{y}_T^i - y_T^i \right)^2$$

□

- 논문 수식 → 코드

$$e_t^k = v_e^+ \tanh(W_e[lh_{t-1}; s_{t-1}] + \text{Ш}_e x^k])  \qquad \text{eqn 8}$$

$$\alpha_t^k = \frac{\exp(e_t^k)}{\sum_{l=1}^{n} \exp(e_t^i)} \quad : \text{attention weight} \qquad \text{eqn 9}$$
$$\text{(softmax function applied)}$$

$lh_{t-1}$          (batch size, m) at time t-1

$s_{t-1}$          (batch size, m) at time t-1

$[lh_{t-1}; s_{t-1}]$      (batch size, 2m)

K. repeat / tf. concat    (batch size, n, 2m)

$W_e[lh_{t-1}; s_{t-1}]$     (batch size, n, T)

  ↳ Dense(T)

permute

$x^k$             (batch size, T, n)

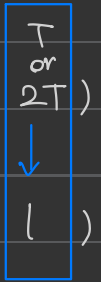$\text{Ш}_e x^k$          (batch size, n, T)

  ↳ Dense(T)

① add($W_e[lh_{t-1}; s_{t-1}]$, $\text{Ш}_e x^k$)     (batch size, n, T)

② concat($W_e[lh_{t-1}; s_{t-1}]$, $\text{Ш}_e x^k$)   (batch size, n, 2T)

$tanh \ (W_e [lh_{t-1} ; S_{t-1}] + U_e x^k ])$

(batch size, n, T or 2T )

$\overline{v_e^+} \ tanh \ (W_e [lh_{t-1} ; S_{t-1}] + U_e x^k ])$

↳ Dense(1)

(batch size, n, 1 )

permute

$\alpha_t^k = \dfrac{exp(e_t^k)}{\sum_{i=1}^{n} exp(e_t^i)}$     (batch size, 1, n)

$$\widetilde{X}_t = (\alpha_t^1 x_t^1, \; \alpha_t^2 x_t^2, \; \cdots, \; \alpha_t^n x_t^n)^t \qquad \text{eqn 10}$$

$$h_t = f_1(h_{t-1}, \; \widetilde{X}_t) \qquad \text{eqn 11}$$