

We can use a Timestep of the Data, of length T , to predict the next K values in the dataset

$$use \vec{x}^t = \begin{bmatrix} x_t \\ x_{t+1} \\ \dots \\ x_{t+T-1} \end{bmatrix} \text{ to predict } \hat{y}^t = \begin{bmatrix} \tilde{x}_{t+T} \\ \tilde{x}_{t+T+1} \\ \dots \\ \tilde{x}_{t+T+k-1} \end{bmatrix}$$

In some way, this “memory element” should also store some information about the prediction \tilde{y}^t . So we can compute \tilde{y}^t as some linear combination of \vec{h}^t .

Diagram illustrating the sequential nature of an RNN. It shows two time steps, t and $t+1$.

At time t :

- Input \vec{x}^t is processed by function $f(\vec{x}^t, \vec{h}^t)$ to produce hidden state \vec{h}^t .
- \vec{h}^t is then combined linearly (Some Linear combination) to produce output \tilde{y}^t .

At time $t+1$:

- Input \vec{x}^{t+1} is processed by function $f(\vec{x}^{t+1}, \vec{h}^t)$, where the hidden state from the previous step (\vec{h}^t) is used.
- This produces \vec{h}^{t+1} , which is then combined linearly (Some Linear combination) to produce output \tilde{y}^{t+1} .

The diagram shows that the hidden state \vec{h}^t is passed to the next time step's function f , demonstrating the sequential dependency.

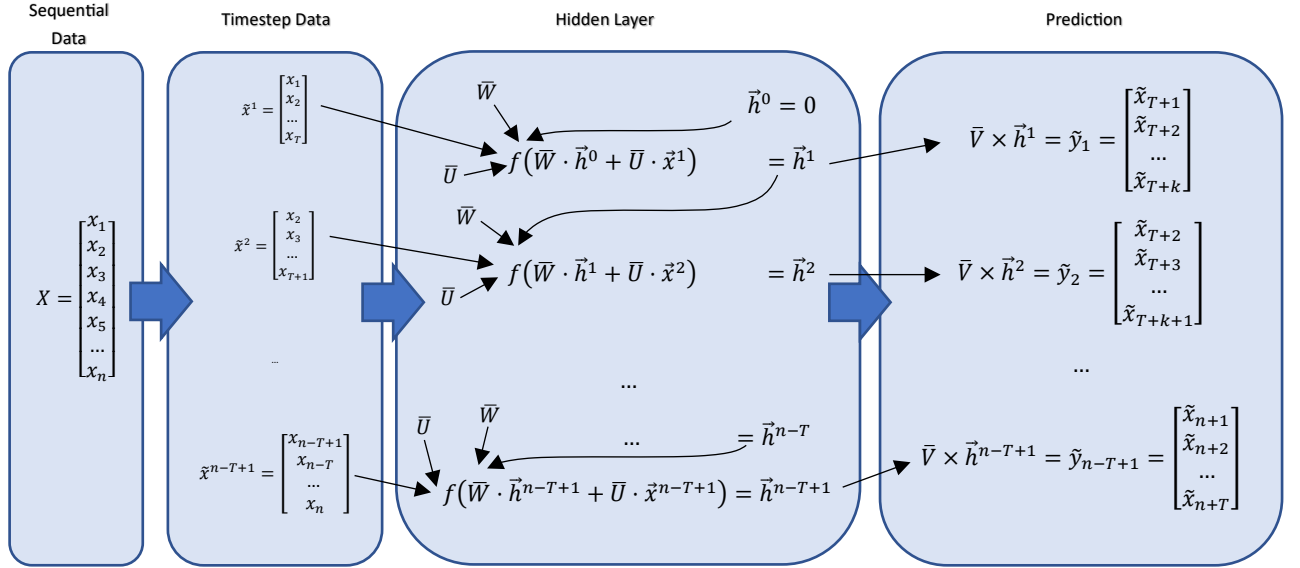
$$\vec{h}^t = f(\bar{W}\vec{h}^{t-1} + \bar{U}\vec{x}^t), \quad \tilde{y} = \bar{V}\vec{h}^t$$

- $\bar{W}_{J \times J} \Rightarrow$ Weighing Parameter for the previous hidden vector \vec{h}^{t-1}
- $\bar{U}_{J \times T} \Rightarrow$ Weighing Patameter for the timestep data \vec{x}^t
- $\bar{V}_{K \times J} \Rightarrow$ Weighing Paramter for the hidden vector \vec{h}^t to compute the prediction \tilde{y}^t
- $\vec{x}^t \in R^T \Rightarrow$ Timestep Data
- $\vec{h}^t \in R^J \Rightarrow$ Hidden Vector
- $\tilde{y}^t \in R^K \Rightarrow$ Prediction

- $T \Rightarrow$ Length of Timestep
- $K \Rightarrow$ Output Dimension
- $J \Rightarrow$ Dimension of hidden vector

Activation function (use sigma function)

$$f(x) = \sigma(x) = \frac{1}{1 + \exp(-x)}$$



Based on the prediction \tilde{y}^t , we can compute the loss function from the actual values of

$$\vec{y}_i = \begin{bmatrix} x_{t+T} \\ x_{t+T+2} \\ \dots \\ x_{t+T+k-1} \end{bmatrix} \text{ as follows:}$$

$$L = \frac{1}{2} \sum_{t=1}^{n-T+1} (\tilde{y}_t - y_t)^2$$

In finding the partial derivatives, it is worth noting the derivative of the sigma function can be simplified as:

$$f'(x) = \sigma'(x) = \sigma(x)(1 - \sigma(x)) = \vec{h}^t(1 - \vec{h}^t)$$

From the loss function, we can thus compute the following partial derivatives of the weighting parameters, $\frac{\partial L}{\partial \vec{W}}$, $\frac{\partial L}{\partial \vec{U}}$, $\frac{\partial L}{\partial \vec{V}}$:

$$\frac{dL_t}{dV_{\alpha\beta}} = \frac{\partial L_t}{\partial \tilde{y}_t} \frac{\partial \tilde{y}_t}{\partial V_{\alpha\beta}} = (\tilde{y}_t - y_t) h_k$$

$$\frac{dL}{d\vec{V}} = \sum_{t=1}^{n-T+1} (\tilde{y}_t - y_t) h_k$$

$$\begin{aligned}
\frac{dL_i}{dU_{\alpha\beta}} &= \frac{\partial L_i}{\partial \tilde{y}_j} \frac{\partial \tilde{y}_j}{\partial h_k} \frac{\partial h_k}{\partial U_{\alpha\beta}} \\
&= (\tilde{y}_i - y_i)(V_{ij}) \left(f'(\bar{W}\vec{h}^{t-1} + \bar{U}\vec{x}^t) \right) (\vec{x}^i) \\
&= (\tilde{y}_i - y_i)(V_{ij}) \left(\vec{h}^t(1 - \vec{h}^t) \right) (\vec{x}^i) \\
\frac{dL}{dU} &= \sum_{i=1}^{n-T+1} (\tilde{y}_i - y_i)(V_{ij}) \left(\vec{h}^t(1 - \vec{h}^t) \right) (\vec{x}^i) \\
\frac{dL_i}{dW_{\alpha\beta}} &= \frac{\partial L_i}{\partial \tilde{y}_j} \frac{\partial \tilde{y}_j}{\partial h_k} \frac{\partial h_k}{\partial W_{\alpha\beta}} \\
&= (\tilde{y}_i - y_i)(V_{ij}) \left(f'(\bar{W}\vec{h}^{t-1} + \bar{U}\vec{x}^t) \right) (\vec{h}^{i-1}) \\
&= (\tilde{y}_i - y_i)(V_{ij}) \left(\vec{h}^t(1 - \vec{h}^t) \right) (\vec{h}^{i-1}) \\
\frac{dL}{dW} &= \sum_{i=1}^{n-T+1} (\tilde{y}_i - y_i)(V_{ij}) \left(\vec{h}^t(1 - \vec{h}^t) \right) (\vec{h}^{i-1})
\end{aligned}$$

Gradient Descent

$$\begin{aligned}
W_{\alpha\beta} &\rightarrow W_{\alpha\beta} - \varepsilon \frac{dL}{dW_{\alpha\beta}} \\
U_{\alpha\beta} &\rightarrow U_{\alpha\beta} - \varepsilon \frac{dL}{dU_{\alpha\beta}} \\
V_{\alpha\beta} &\rightarrow V_{\alpha\beta} - \varepsilon \frac{dL}{dV_{\alpha\beta}}
\end{aligned}$$

Algorithm

Initialize Parameters

- $\bar{W} = 0, \bar{V} = 0, \bar{U} = 0$
- $\vec{h}^0 = 0$

Iterate for N epochs

For every n-T+1 combination of $\vec{x}^t = \begin{bmatrix} x_t \\ x_{t+1} \\ \dots \\ x_{t+T} \end{bmatrix},$

Find hidden vector from previous hidden vector

$$\vec{h}^t = f(\bar{W}\vec{h}^{t-1} + \bar{U}\vec{x}^t)$$

Compute predictions $\tilde{y}_t = \bar{V}\vec{h}^t = \bar{V}f(\bar{W}\vec{h}^{t-1} + \bar{U}\vec{x}^t) \forall t = 1, 2, \dots, n - T + 1$

Compute Partial Derivatives and update Parameters

$$\frac{dL}{d\bar{V}} = \sum_{i=1}^{n-T+1} (\tilde{y}_i - y_i) h_k$$

$$\bar{V} \rightarrow \bar{V} - \varepsilon \frac{dL}{d\bar{V}}$$

$$\frac{dL}{d\bar{U}} = \sum_{i=1}^{n-T+1} (\tilde{y}_t - y_t) (V_{ij}) (\vec{h}^t (1 - \vec{h}^t)) (\vec{x}^t)$$

$$\bar{U} \rightarrow \bar{U} - \varepsilon \frac{dL}{d\bar{U}}$$

$$\frac{dL}{d\bar{W}} = \sum_{i=1}^{n-T+1} (\tilde{y}_i - y_i) (V_{ij}) (\vec{h}^t (1 - \vec{h}^t)) (\vec{h}^{i-1})$$

$$\bar{W} \rightarrow \bar{W} - \varepsilon \frac{dL}{d\bar{W}}$$