

Lecture notes on recurrent neural networks (RNNs)

1.1 Simple RNNs

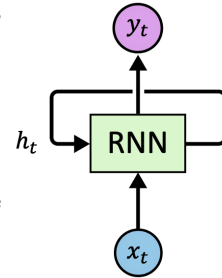
Mathematical form

Assume we want to build a neural network to process a data sequence $\mathbf{x}_{1:T} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$. Recurrent neural networks (RNNs) are neural networks suited for processing sequential data, which, if well trained, can model dependencies within a sequence of arbitrary length.

We also assume a supervised learning task which aims to learn the mapping from inputs $\mathbf{x}_{1:T}$ to outputs $\mathbf{y}_{1:T} = (\mathbf{y}_1, \dots, \mathbf{y}_T)$. Then a simple RNN computes the following mapping for $t = 1, \dots, T$:

$$\mathbf{h}_t = \phi_h(W_h \mathbf{h}_{t-1} + W_x \mathbf{x}_t + \mathbf{b}_h), \quad (1)$$

$$\mathbf{y}_t = \phi_y(W_y \mathbf{h}_t + \mathbf{b}_y). \quad (2)$$



Here the network parameters are $\theta = \{W_h, W_x, W_y, \mathbf{b}_h, \mathbf{b}_y\}$, ϕ_h and ϕ_y are the non-linear activation functions for the hidden state \mathbf{h}_t and the output y_t , respectively. For $t = 1$ the convention is to set $\mathbf{h}_0 = \mathbf{0}$ so that $\mathbf{h}_1 = \phi_h(W_x \mathbf{x}_1 + \mathbf{b}_h)$; alternatively \mathbf{h}_0 can also be added to θ as a learnable parameter.

Back-propagation through time (BPTT)

A loss function is required for training an RNN. Assuming the following loss function to minimise:

$$\mathcal{L}(\theta) = \sum_{t=1}^T \mathcal{L}(\mathbf{y}_t). \quad (3)$$

This is a common form for the loss function in many sequential modelling tasks such as video/audio sequence reconstruction. The derivative of the loss function w.r.t. θ is $\frac{d}{d\theta} \mathcal{L}(\theta) = \sum_{t=1}^T \frac{d}{d\theta} \mathcal{L}(\mathbf{y}_t)$, therefore it remains to compute $\frac{d}{d\theta} \mathcal{L}(\mathbf{y}_t)$ for $\theta = \{W_h, W_x, W_y, \mathbf{b}_h, \mathbf{b}_y\}$.

- Derivative of $\mathcal{L}(\mathbf{y}_t)$ w.r.t. W_y and \mathbf{b}_y :

$$\frac{d\mathcal{L}(\mathbf{y}_t)}{dW_y} = \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{y}_t} \frac{d\mathbf{y}_t}{dW_y}, \quad \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{b}_y} = \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{y}_t} \frac{d\mathbf{y}_t}{d\mathbf{b}_y}.$$

- Derivative of $\mathcal{L}(\mathbf{y}_t)$ w.r.t. W_x and \mathbf{b}_h :

$$\frac{d\mathcal{L}(\mathbf{y}_t)}{dW_x} = \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{y}_t} \frac{d\mathbf{y}_t}{d\mathbf{h}_t} \frac{d\mathbf{h}_t}{dW_x}, \quad \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{b}_h} = \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{y}_t} \frac{d\mathbf{y}_t}{d\mathbf{h}_t} \frac{d\mathbf{h}_t}{d\mathbf{b}_h}.$$

Here W_x and \mathbf{b}_h contributes to \mathbf{h}_t in two ways: both direct and indirect contributions, where the latter is through \mathbf{h}_{t-1} . This means:

$$\frac{d\mathbf{h}_t}{dW_x} = \frac{\partial \mathbf{h}_t}{\partial W_x} + \frac{d\mathbf{h}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{h}_{t-1}}{dW_x}, \quad \frac{d\mathbf{h}_t}{d\mathbf{b}_h} = \frac{\partial \mathbf{h}_t}{\partial \mathbf{b}_h} + \frac{d\mathbf{h}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{h}_{t-1}}{d\mathbf{b}_h}.$$

Derivations of these derivatives requires the usage of chain rule which will be explained next.

- Derivative of $\mathcal{L}(\mathbf{y}_t)$ w.r.t. W_h : by chain rule, we have

$$\frac{d\mathcal{L}(\mathbf{y}_t)}{dW_h} = \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{h}_t} \frac{d\mathbf{h}_t}{dW_h}$$

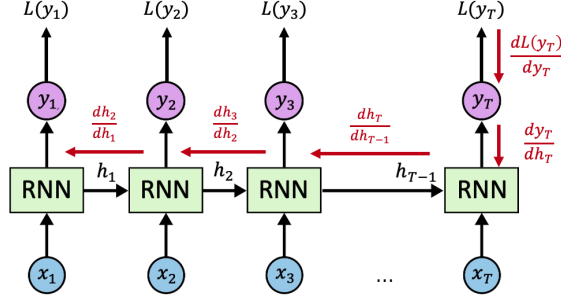


Figure 1: Visualising Back-propagation through time (BPTT) without truncation. The black arrows show forward pass computations, while the red arrows show the gradient back-propagation in order to compute $\nabla_{W_h} \mathcal{L}(y_t)$.

where $\frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{h}_t} = \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{y}_t} \frac{d\mathbf{y}_t}{d\mathbf{h}_t}$. Importantly, here the entries in the Jacobian $\frac{d\mathbf{h}_t}{dW_h}$ contains the *total gradient* of $\mathbf{h}_t[i]$ w.r.t. $W_h[m, n]$. It remains to compute $\frac{d\mathbf{h}_t}{dW_h}$ and notice that \mathbf{h}_t depends on \mathbf{h}_{t-1} which also depends on W_h :

$$\frac{d\mathbf{h}_t}{dW_h} = \frac{\partial \mathbf{h}_t}{\partial W_h} + \frac{d\mathbf{h}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{h}_{t-1}}{dW_h}. \quad (4)$$

Here the entries in $\frac{\partial \mathbf{h}_t}{\partial W_h}$ contains *partial gradient* only (by treating \mathbf{h}_{t-1} as a constant w.r.t. W_h , note that \mathbf{h}_t depends on \mathbf{h}_{t-1}). By expanding the $\frac{d\mathbf{h}_{t-1}}{dW_h}$ term further, we have:

$$\begin{aligned} \frac{d\mathbf{h}_t}{dW_h} &= \frac{\partial \mathbf{h}_t}{\partial W_h} + \frac{d\mathbf{h}_t}{d\mathbf{h}_{t-1}} \frac{\partial \mathbf{h}_{t-1}}{\partial W_h} + \frac{d\mathbf{h}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{h}_{t-1}}{d\mathbf{h}_{t-2}} \frac{d\mathbf{h}_{t-2}}{dW_h} = \dots \\ &= \sum_{\tau=1}^t \left(\prod_{l=\tau}^{t-1} \frac{d\mathbf{h}_{l+1}}{d\mathbf{h}_l} \right) \frac{\partial \mathbf{h}_\tau}{\partial W_h}, \end{aligned} \quad (5)$$

with the convention that when $\tau = t$, $\prod_{l=t}^{t-1} \frac{d\mathbf{h}_{l+1}}{d\mathbf{h}_l} = 1$. This means the chain rule of the gradients needs to be computed in an reversed order from time $t = T$ to time $t = 1$, hence the name *Back-propagation through time (BPTT)*. A visualisation of BPTT is provided in Figure 1. Truncation with length L might be applied to this back-propagation procedure, and with *truncated BPTT* the gradient is computed as

$$\text{truncate}\left[\frac{d\mathbf{h}_t}{dW_h}\right] = \sum_{\tau=\max(1, t-L)}^t \left(\prod_{l=\tau}^{t-1} \frac{d\mathbf{h}_{l+1}}{d\mathbf{h}_l} \right) \frac{\partial \mathbf{h}_\tau}{\partial W_h}.$$

Gradient vanishing/explosion issues

Simple RNNs are often said to suffer from gradient vanishing or gradient explosion issues. To understand this, notice that

$$\frac{d\mathbf{h}_{l+1}}{d\mathbf{h}_l}^\top = \phi'_h(W_h \mathbf{h}_l + W_x \mathbf{x}_{l+1} + \mathbf{b}_h) \odot W_h, \quad (6)$$

Here $\phi'_h(W_h \mathbf{h}_l + W_x \mathbf{x}_{l+1} + \mathbf{b}_h)$ denotes a vector containing element-wise derivatives, and we reload the element-wise product operator for vector $\mathbf{a} \in \mathbb{R}^{d \times 1}$ and $\mathbf{B} \in \mathbb{R}^{d \times d'}$ as the “broadcasting element-wise product” $\mathbf{a} \odot \mathbf{B} := [\underbrace{\mathbf{a}, \dots, \mathbf{a}}_{\text{repeat } d' \text{ times}}] \odot \mathbf{B}$. This means $\prod_{l=\tau}^{t-1} \frac{d\mathbf{h}_{l+1}}{d\mathbf{h}_l}$ contains products of $t - \tau$ copies of W_h and the derivative $\phi'_h(\cdot)$ at time steps $l = \tau, \dots, t - 1$.

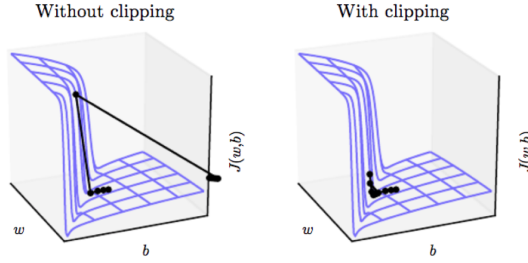


Figure 2: Visualising the gradient step with/without gradient clipping. Source: [Goodfellow et al. \[2016\]](#).

Now consider a simple case where $\phi_h(\cdot)$ is an identity mapping so that $\phi'_h(\cdot) = 1$. We further assume the hidden states have scalar values, i.e. $\dim(\mathbf{h}_t) = 1$. Then we have $\prod_{l=\tau}^{t-1} \frac{dh_{l+1}}{dh_l} = (W_h^{t-\tau})^\top$ which can vanish or explode when $t - \tau$ is large, depending on whether $W_h < 1$ or not. In the general case when W_h is a matrix, depending on whether the largest singular value (i.e. maximum of the absolute values of the largest and smallest eigenvalues) of W_h is smaller or larger than 1, the spectral norm of $\prod_{l=\tau}^{t-1} \frac{dh_{l+1}}{dh_l} = (W_h^{t-\tau})^\top$ will vanish or explode when $t - \tau$ increases. When $\phi_h(\cdot)$ is selected as the sigmoid function or the hyperbolic tangent function, gradient vanishing problem can still happen. Take hyperbolic tangent function as an example. When entries in \mathbf{h}_t is close to ± 1 , then $\phi'_h(\cdot) \approx 0$, i.e. the derivative is saturated. Multiplying several of such saturated derivatives together also leads to the gradient vanishing problem.

Some tricks to fix the gradient vanishing/explosion issues

There are a handful of empirical tricks to fix the gradient vanishing/explosion issues discussed above.

- Gradient clipping:

This trick is often used to prevent the gradient from explosion. With a fixed hyper-parameter γ , a gradient \mathbf{g} is clipped when $\|\mathbf{g}\| > \gamma$:

$$\mathbf{g} \leftarrow \frac{\gamma}{\|\mathbf{g}\|} \mathbf{g}.$$

This trick ensures the gradients used in optimisation has their maximum norm bounded by a pre-defined hyper-parameter. It introduces biases in the gradient-based optimisation procedure, but in certain cases it can be beneficial. Figure 2 visualises such an example, where with gradient clipping, the updates can stay in the valley of the loss function.

- Good initialisation of the recurrent weight matrix W_h :
The IRNN approach [[Le et al., 2015](#)] uses ReLU activation's for ϕ_h and initialise $W_h = \mathbf{I}$, $\mathbf{b}_h = \mathbf{0}$. This makes $\phi'_h(t) = \delta(t > 0)$ and $\frac{dh_{l+1}}{dh_l} = \delta(W_x \mathbf{x}_{l+1} > 0)$ at initialisation. While there is no guarantee of eliminating the gradient vanishing/explosion problem during the whole course of training, empirically RNNs with this trick have achieved competitive performance to LSTMs in a variety of tasks.
- Alternatively, one can construct the recurrent weight matrix W_h to be orthogonal or unitary matrix. See e.g. [Saxe et al. \[2014\]](#); [Arjovsky et al. \[2016\]](#) for examples.

1.2 Long Short-Term Memory (LSTM)

Mathematical form

The Long Short-term Memory [[Hochreiter and Schmidhuber, 1997](#)] was proposed with the motivation of addressing the gradient vanishing/explosion problem. It introduces *memory cell states* and *gates*

to control the error flows, in detail the computation goes as follows (with $\sigma(\cdot)$ as sigmoid function):

$$\mathbf{f}_t : \text{forget gate} \quad \mathbf{f}_t = \sigma(W_f \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f) \quad (7)$$

$$\mathbf{i}_t : \text{input gate} \quad \mathbf{i}_t = \sigma(W_i \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i) \quad (8)$$

$$\mathbf{o}_t : \text{output gate} \quad \mathbf{o}_t = \sigma(W_o \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o) \quad (9)$$

$$\mathbf{x}_t : \text{input} \quad \tilde{\mathbf{c}}_t = \tanh(W_c \cdot [\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_c) \quad (10)$$

$$\mathbf{c}_t : \text{memory cell state} \quad \mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t \quad (11)$$

$$\mathbf{h}_t : \text{hidden state} \quad \mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t) \quad (12)$$

The parameters of an LSTM are therefore $\boldsymbol{\theta} = \{W_f, W_i, W_o, W_c, \mathbf{b}_f, \mathbf{b}_i, \mathbf{b}_o, \mathbf{b}_c\}$. Again by convention, \mathbf{h}_0 and \mathbf{c}_0 are either set to zero vectors or added to the learnable parameters. We note that if the initial cell state \mathbf{c}_0 is initialised to zero, then we have the elements in \mathbf{c}_t and \mathbf{h}_t bounded within $(-1, 1)$. The output y_t can be produced by a 1-layer neural network similar to the simple RNN case: $\mathbf{y}_t = \phi_y(W_y \mathbf{h}_t + \mathbf{b}_y)$.

*Gradient computation

Readers are encouraged to derive themselves the gradient of $\mathcal{L}(\mathbf{y}_t)$ with respect to $\boldsymbol{\theta}$. Specifically for the recurrent weight matrix W_c , computing the derivative $\frac{d\mathcal{L}(\mathbf{y}_t)}{dW_c}$ requires the following terms:

$$\frac{d\mathcal{L}(\mathbf{y}_t)}{dW_c} = \frac{d\mathcal{L}(\mathbf{y}_t)}{d\mathbf{h}_t} \frac{d\mathbf{h}_t}{dW_c} \quad (13)$$

$$\frac{d\mathbf{h}_t}{dW_c} = \mathbf{o}_t \odot \frac{d\tanh(\mathbf{c}_t)}{dW_c} + \tanh(\mathbf{c}_t) \odot \frac{d\mathbf{o}_t}{dW_c} \quad (14)$$

$$\frac{d\mathbf{o}_t}{dW_c} = \frac{d\mathbf{o}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{h}_{t-1}}{dW_c} \quad (15)$$

$$\frac{d\mathbf{c}_t}{dW_c} = \mathbf{f}_t \odot \frac{d\mathbf{c}_{t-1}}{dW_c} + \mathbf{c}_{t-1} \odot \frac{d\mathbf{f}_t}{dW_c} + \mathbf{i}_t \odot \frac{d\tilde{\mathbf{c}}_t}{dW_c} + \tilde{\mathbf{c}}_t \odot \frac{d\mathbf{i}_t}{dW_c}. \quad (16)$$

As \mathbf{h}_{t-1} also depends on \mathbf{c}_{t-1} and \mathbf{o}_{t-1} , it means that

$$\frac{d\mathbf{f}_t}{dW_c} = \frac{d\mathbf{f}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{h}_{t-1}}{dW_c} = \frac{d\mathbf{f}_t}{d\mathbf{h}_{t-1}} \left(\mathbf{o}_{t-1} \odot \frac{d\tanh(\mathbf{c}_{t-1})}{dW_c} + \tanh(\mathbf{c}_{t-1}) \odot \frac{d\mathbf{o}_{t-1}}{dW_c} \right) \quad (17)$$

$$\frac{d\mathbf{i}_t}{dW_c} = \frac{d\mathbf{i}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{h}_{t-1}}{dW_c} = \frac{d\mathbf{i}_t}{d\mathbf{h}_{t-1}} \left(\mathbf{o}_{t-1} \odot \frac{d\tanh(\mathbf{c}_{t-1})}{dW_c} + \tanh(\mathbf{c}_{t-1}) \odot \frac{d\mathbf{o}_{t-1}}{dW_c} \right) \quad (18)$$

$$\frac{d\tilde{\mathbf{c}}_t}{dW_c} = \frac{\partial \tilde{\mathbf{c}}_t}{\partial W_c} + \frac{d\tilde{\mathbf{c}}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{h}_{t-1}}{dW_c} = \frac{\partial \tilde{\mathbf{c}}_t}{\partial W_c} + \frac{d\tilde{\mathbf{c}}_t}{d\mathbf{h}_{t-1}} \left(\mathbf{o}_{t-1} \odot \frac{d\tanh(\mathbf{c}_{t-1})}{dW_c} + \tanh(\mathbf{c}_{t-1}) \odot \frac{d\mathbf{o}_{t-1}}{dW_c} \right). \quad (19)$$

Note again the difference between $\frac{d\tilde{\mathbf{c}}_t}{dW_c}$ and $\frac{\partial \tilde{\mathbf{c}}_t}{\partial W_c}$. The former Jacobian $\frac{d\tilde{\mathbf{c}}_t}{dW_c}$ has its entries as the *total gradient* of $\tilde{\mathbf{c}}_t[i]$ w.r.t. $W_c[m, n]$, while the latter partial gradient $\frac{\partial \tilde{\mathbf{c}}_t}{\partial W_c}$ has its entries as the *partial gradient* of $\tilde{\mathbf{c}}_t[i]$ w.r.t. $W_c[m, n]$ (by treating \mathbf{h}_{t-1} as a constant w.r.t. W_c , note that $\tilde{\mathbf{c}}_t$ depends on \mathbf{h}_{t-1} as well). Combining the derivations, we have (notice that $\frac{d\tanh(\mathbf{c}_{t-1})}{dW_c} = \frac{d\tanh(\mathbf{c}_{t-1})}{d\mathbf{c}_{t-1}} \frac{d\mathbf{c}_{t-1}}{dW_c}$):

$$\frac{d\mathbf{o}_t}{dW_c} = \frac{d\mathbf{o}_t}{d\mathbf{h}_{t-1}} \left(\mathbf{o}_{t-1} \odot \frac{d\tanh(\mathbf{c}_{t-1})}{dW_c} + \tanh(\mathbf{c}_{t-1}) \odot \frac{d\mathbf{o}_{t-1}}{dW_c} \right) \quad (20)$$

$$\frac{d\mathbf{c}_t}{dW_c} = \underbrace{\left(\mathbf{f}_t + \mathbf{o}_{t-1} \odot \frac{d\tanh(\mathbf{c}_{t-1})}{d\mathbf{c}_{t-1}} \odot \frac{d\mathbf{c}_{t-1}}{d\mathbf{h}_{t-1}} \right)}_{=\frac{d\mathbf{c}_t}{d\mathbf{c}_{t-1}}} \frac{d\mathbf{c}_{t-1}}{dW_c} + \tanh(\mathbf{c}_{t-1}) \odot \frac{d\mathbf{c}_t}{d\mathbf{h}_{t-1}} \frac{d\mathbf{o}_{t-1}}{dW_c} + \mathbf{i}_t \odot \frac{\partial \tilde{\mathbf{c}}_t}{\partial W_c} \quad (21)$$

$$\frac{d\mathbf{c}_t}{d\mathbf{h}_{t-1}} = \mathbf{c}_{t-1} \odot \frac{d\mathbf{f}_t}{d\mathbf{h}_{t-1}} + \tilde{\mathbf{c}}_t \odot \frac{d\mathbf{i}_t}{d\mathbf{h}_{t-1}} + \mathbf{i}_t \odot \frac{d\tilde{\mathbf{c}}_t}{d\mathbf{h}_{t-1}}. \quad (22)$$

This means for computing $\frac{dc_t}{dW_c}$ it requires computing

$$\prod_{l=\tau}^{t-1} \frac{dc_{l+1}}{dc_l} = \prod_{l=\tau}^{t-1} [f_{l+1} + o_l \odot \frac{d \tanh(c_l)}{dc_l} \odot \frac{dc_{l+1}}{dh_l}]$$

for all $\tau = 1, \dots, t$. There is no guarantee that this term will not vanish or explode, however the usage of forget gates makes the issue less severe. To see this, notice that by expanding the product term above, we have that it contains terms proportional to $f_{i+1} \odot \prod_{l=\tau}^i o_l \odot \frac{dc_{l+1}}{dh_l}$ for $i = \tau + 1, \dots, t - 1$. So if in the forward pass the network sets $f_{i+1} \rightarrow 0$ (i.e. forgetting the previous cell state), then this will also likely to bring $f_{i+1} \odot \prod_{l=\tau}^i o_l \odot \frac{dc_{l+1}}{dh_l} \approx 0$, which is helpful to cope with the gradient explosion problem. On the other hand, $\frac{dc_t}{dW_c}$ also contains terms proportional to $\prod_{l=\tau+1}^i f_l \odot o_\tau \odot \frac{dc_{\tau+1}}{dh_\tau}$ for $i = \tau + 1, \dots, t - 1$. This means if the network sets $f_l \rightarrow 1$ for $l = \tau + 1, \dots, i$ (i.e. maintaining the cell state until at least time $t = i$), then it will be likely that $\prod_{l=\tau+1}^i f_l \odot o_\tau \odot \frac{dc_{\tau+1}}{dh_\tau} \approx o_\tau \odot \frac{dc_{\tau+1}}{dh_\tau}$, which helps prevent the gradient info at time τ from vanishing when $o_\tau \rightarrow 1$, and thus be helpful to learn longer term dependencies. The gradients $\frac{dc_t}{dW_c}$ and $\frac{do_t}{dW_c}$ also require computing products of $\frac{do_{i+1}}{dh_i}$ and $\frac{dc_{i+1}}{dh_i}$ terms, and analogous analysis can be done for those product terms. It is worth emphasising again that LSTM does NOT solve the gradient vanishing/explosion problem completely, however empirical evidences have shown that it is easier for LSTMs to learn longer term dependencies when compared with the simple RNN.

Gated Recurrent Unit (GRU): a simplified gated RNN

The Gated Recurrent Unit (GRU) [Cho et al., 2014] improves the simple RNN with the gating mechanism as well. Compared with LSTM, GRU removes the input/output gates and the cell state, but still maintains the forgetting mechanism in some form:

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z) \quad (23)$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t] + b_r) \quad (24)$$

$$\tilde{h}_t = \tanh(W_h \cdot [r_t \odot h_{t-1}, x_t] + b_h) \quad (25)$$

$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \tilde{h}_t \quad (26)$$

The network parameters are then $\theta = \{W_z, W_r, W_h, b_z, b_r, b_h\}$. We see here z_t acts as the the update gate which determines the incorporation of the current info to the hidden states, and r_t is named the reset gate which also impact on the maintenance of the historical information.

1.3 Sequence-to-sequence models

Given dataset of input-output sequence pairs $(\mathbf{x}_{1:T}, \mathbf{y}_{1:L})$, the goal of sequence prediction is to build a model $p_\theta(\mathbf{y}_{1:L} | \mathbf{x}_{1:T})$ to fit the data. Note here that the \mathbf{x}, \mathbf{y} sequences might have different lengths, and the input/output length T and L can vary across input-output pairs. So to handle sequence outputs of arbitrary length, we define an *auto-regressive model*

$$p_\theta(\mathbf{y}_{1:L} | \mathbf{x}_{1:T}) = \prod_{l=1}^L p_\theta(\mathbf{y}_l | \mathbf{y}_{<l}, \mathbf{v}), \quad \mathbf{v} = enc(\mathbf{x}_{1:T}). \quad (27)$$

Here $p_\theta(\mathbf{y}_l | \mathbf{y}_{<l}, \mathbf{v})$ is defined by a sequence decoder, e.g. an LSTM:

$$p_\theta(\mathbf{y}_l | \mathbf{y}_{<l}, \mathbf{v}) = p_\theta(\mathbf{y}_l | \mathbf{h}_l^d, \mathbf{c}_l^d), \quad \mathbf{h}_l^d, \mathbf{c}_l^d = LSTM_\theta^{dec}(\mathbf{y}_{<l}), \quad (28)$$

and the decoder LSTM has its internal recurrent states $\mathbf{h}_0^d, \mathbf{c}_0^d$ initialised using the input sequence representation $\mathbf{v} = enc(\mathbf{x}_{1:T})$. The encoder also uses an LSTM, meaning that

$$\mathbf{v} = enc(\mathbf{x}_{1:T}) = NN_\theta(\mathbf{h}_T^e, \mathbf{c}_T^e), \quad \mathbf{h}_T^e, \mathbf{c}_T^e = LSTM_\theta^{enc}(\mathbf{x}_{1:T}). \quad (29)$$

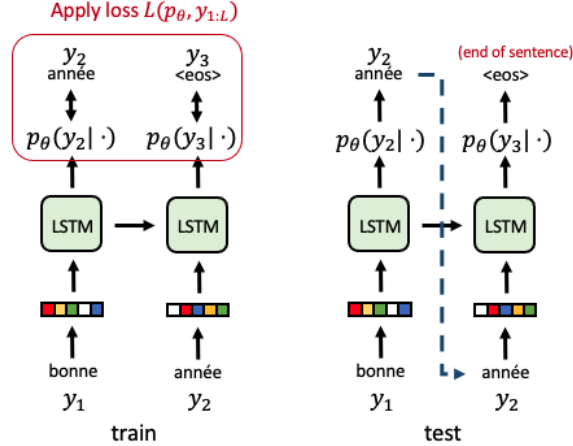


Figure 3: Visualising the forward pass of Seq2Seq model in train (left) and test (right) time.

Regarding the prediction for the first output \mathbf{y}_1 , either $p_{\theta}(\mathbf{y}_1|\mathbf{x}_{1:T})$ can be produced using the last recurrent states of the encoder LSTM (i.e. $p_{\theta}(\mathbf{y}_1|\mathbf{x}_{1:T}) = p_{\theta}(\mathbf{y}_1|\mathbf{h}_T^e, \mathbf{c}_T^e)$), or we can add in a “start of sentence” token as \mathbf{y}_0 and compute the probability vector for \mathbf{y}_1 using the LSTM decoder: $p_{\theta}(\mathbf{y}_1|\mathbf{x}_{1:T}) = p_{\theta}(\mathbf{y}_1|\mathbf{y}_0, \mathbf{v})$. This model is named *Sequence-to-sequence* model or *Seq2Seq* model in short, which is proposed by Sutskever et al. [2014].

The decoder LSTM forward pass in training and test times are different, which is visualised in Figure 3. In training, since maximum likelihood training requires evaluating $p_{\theta}(\mathbf{y}_l|\mathbf{y}_{<l}, \mathbf{v})$ for the output sequences $\mathbf{y}_{1:L}$ from the dataset, this means the inputs to the decoder LSTM are words in the data label sequence, and the output of the LSTM is the probability vector for the current word, which is then used to compute the MLE objective (i.e. negative cross-entropy). In test time, however, there is no ground truth output sequence provided, therefore the input to the decoder LSTM at step l is the *predicted* word $\mathbf{y}_{l-1} \sim p_{\theta}(\mathbf{y}_{l-1}|\mathbf{y}_{<l-1}, \mathbf{v})$. In practice prediction is done by e.g. beam search rather than naive sequential sampling [Sutskever et al., 2014].

In NLP applications such as machine translation, both $\mathbf{x}_{1:T}$ and $\mathbf{y}_{1:L}$ are sequence of words which cannot be directly processed by neural networks. Instead each \mathbf{x}_t (\mathbf{y}_l) needs to be mapped to a real-value vector before feeding it to the encoder (decoder) LSTM. A naive approach is to use one-hot encoding: assume that the input sentence is in English and we have a sorted English vocabulary of size V , then

$$\mathbf{x}_t \rightarrow \underbrace{(0, 0, \dots, 1, \dots, 0)}_{k-1}, \quad \text{if } \mathbf{x}_t \text{ is the } k\text{th word in the vocabulary.}$$

This is clearly inefficient since English vocabulary has tens of thousands of words. Instead, it is recommended to map the words to their *word embeddings* using word2vec [Mikolov et al., 2013a] or GloVe [Pennington et al., 2014], which has much smaller dimensions. But more importantly, semantics are preserved to some extent in these word embeddings, e.g. it has been shown that vector calculus results like “emb(king) - emb(male) + emb(female) = emb(queen)” hold for word2vec embeddings.

In many NLP applications the decoder output is a probability vector $p_{\theta}(\mathbf{y}_l|\mathbf{y}_{<l}, \mathbf{v})$ which specifies the predictive probability of each of the words in the vocabulary. When the vocabulary is large (which is often so), applying softmax to obtain the probability vector can be very challenging. Interested readers can check e.g. hierarchical softmax [Mikolov et al., 2013a] and negative sampling [Mikolov et al., 2013b] for solutions to mitigate this issue.

1.4 *Generative models for sequences

To generate sequential data such as video, text and audio, one needs to build a generative model $p_{\theta}(\mathbf{x}_{1:T})$ and train it with e.g. (approximate) maximum likelihood. In the following we discuss two types of latent variable models that are often used in sequence generation tasks.

Sequence VAE with global latent variables

Similar to VAEs for image generation, one can define a latent variable model with a global latent variable for sequence generation [Fabius and van Amersfoort, 2014; Bowman et al., 2016]:

$$p_{\theta}(\mathbf{x}_{1:T}) = \int p_{\theta}(\mathbf{x}_{1:T}|\mathbf{z})p(\mathbf{z})d\mathbf{z}. \quad (30)$$

If variational lower-bound is used for training, then this also requires an approximate posterior $q_{\phi}(\mathbf{z}|\mathbf{x}_{1:T})$ to be optimised:

$$\phi^*, \theta^* = \arg \max \mathcal{L}(\phi, \theta), \quad \mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{\text{data}}(\mathbf{x}_{1:T})} \underbrace{[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}_{1:T})}[\log p_{\theta}(\mathbf{x}_{1:T}|\mathbf{z})] - \text{KL}[q_{\phi}(\mathbf{z}|\mathbf{x}_{1:T})||p(\mathbf{z})]]}_{:=\mathcal{L}(\mathbf{x}_{1:T}, \phi, \theta)}. \quad (31)$$

Now it remains to define the encoder and decoder distributions, such that they can process sequence of any length. This can be achieved using e.g. LSTMs to define an auto-regressive decoder:

$$p_{\theta}(\mathbf{x}_{1:T}|\mathbf{z}) = \prod_{t=1}^T p_{\theta}(\mathbf{x}_t|\mathbf{x}_{<t}, \mathbf{z}), \quad p_{\theta}(\mathbf{x}_1|\mathbf{x}_{<1}, \mathbf{z}) = p_{\theta}(\mathbf{x}_1|\mathbf{z}). \quad (32)$$

with the distributional parameters of $p_{\theta}(\mathbf{x}_t|\mathbf{x}_{<t}, \mathbf{z})$ defined by $LSTM_{\theta}(\mathbf{x}_{<t})$ which has its recurrent states $\mathbf{h}_0, \mathbf{c}_0$ initialised using \mathbf{z} . For the encoder, LSTMs can also be used to process the input:

$$q_{\phi}(\mathbf{z}|\mathbf{x}_{1:T}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\phi}(\mathbf{x}_{1:T}), \text{diag}(\boldsymbol{\sigma}_{\phi}^2(\mathbf{x}_{1:T}))), \quad \boldsymbol{\mu}_{\phi}(\mathbf{x}_{1:T}), \log \boldsymbol{\sigma}_{\phi}(\mathbf{x}_{1:T}) = LSTM_{\phi}(\mathbf{x}_{1:T}). \quad (33)$$

State-space models

State-space models assume that for every observation \mathbf{x}_t at time t , there is a latent variable \mathbf{z}_t that generates it, and the sequence dynamic model is defined in the latent space rather than in the observation space. In detail, a *prior dynamic model* is assumed on the transitions of the latent states \mathbf{z}_t , often in an auto-regressive way:

$$p_{\theta}(\mathbf{z}_{1:T}) = \prod_{t=1}^T p_{\theta}(\mathbf{z}_t|\mathbf{z}_{<t}), \quad p_{\theta}(\mathbf{z}_1|\mathbf{z}_{<1}) = p_{\theta}(\mathbf{z}_1). \quad (34)$$

The observation \mathbf{x}_t at time t is assumed to be conditionally dependent on \mathbf{z}_t only, and this conditional distribution is also called the *emission model*:

$$p_{\theta}(\mathbf{x}_t|\mathbf{z}_{1:T}) = p_{\theta}(\mathbf{x}_t|\mathbf{z}_t). \quad (35)$$

Combining both definitions, we have the sequence generative model defined as

$$p_{\theta}(\mathbf{x}_{1:T}) = \int \prod_{t=1}^T p_{\theta}(\mathbf{x}_t|\mathbf{z}_t)p_{\theta}(\mathbf{z}_t|\mathbf{z}_{<t})d\mathbf{z}_{1:T}. \quad (36)$$

A variational lower-bound objective for training this state-space model require an approximate posterior distribution $q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$:

$$\begin{aligned} \mathcal{L}(\phi, \theta) &= \mathbb{E}_{p_{\text{data}}(\mathbf{x}_{1:T})} [\mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})} [\log p_{\theta}(\mathbf{x}_{1:T}|\mathbf{z}_{1:T})] - \text{KL}[q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})||p_{\theta}(\mathbf{z}_{1:T})]] \\ &= \mathbb{E}_{p_{\text{data}}(\mathbf{x}_{1:T})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})} \left[\sum_{t=1}^T \log p_{\theta}(\mathbf{x}_t|\mathbf{z}_t) \right] - \text{KL}[q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})||p_{\theta}(\mathbf{z}_{1:T})] \right]. \end{aligned} \quad (37)$$

Different from the image generation case, now the prior distribution $p_{\theta}(\mathbf{z}_{1:T})$ also has learnable parameters in θ , so in this case it is less appropriate to view this KL term as a “regulariser”.

The expanded expression for the variational lower-bound depends on the definition of the encoder distribution $q_{\theta}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T})$. The simplest solution is to use a factorised approximate posterior

$$q_{\phi}(\mathbf{z}_{1:T}|\mathbf{x}_{1:T}) = \prod_{t=1}^T q(\mathbf{z}_t|\mathbf{x}_{\leq t}), \quad (38)$$

and the variational lower-bound becomes

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{p_{\text{data}}(\mathbf{x}_{1:T})} \left[\sum_{t=1}^T \mathbb{E}_{q(\mathbf{z}_{<t}|\mathbf{x}_{<t})} \left[\underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}_t|\mathbf{x}_{\leq t})} [\log p_{\theta}(\mathbf{x}_t|\mathbf{z}_t)] - \text{KL}[q_{\phi}(\mathbf{z}_t|\mathbf{x}_{\leq t})||p_{\theta}(\mathbf{z}_t|\mathbf{z}_{<t})]}_{:=\mathcal{L}(\mathbf{x}_t, \phi, \theta, \mathbf{z}_{<t})} \right] \right]. \quad (39)$$

We see that the term $\mathcal{L}(\mathbf{x}_t, \phi, \theta, \mathbf{z}_{<t})$ resembles the VAE objective in the image generation case, except that the prior distribution $p_{\theta}(\mathbf{z}_t|\mathbf{z}_{<t})$ is conditioned on the previous latent states $\mathbf{z}_{<t}$ rather than a standard Gaussian, and the q distribution takes $\mathbf{x}_{\leq t}$ as the input rather than a single frame \mathbf{x}_t .

Neural networks can be used to construct the conditional distributions in the following way. The distributional parameters (e.g. mean and variance) of the emission model $p_{\theta}(\mathbf{x}_t|\mathbf{z}_t)$ can be defined by a neural network transformation of \mathbf{z}_t , similar to deep generative models for images. The prior dynamic model $p_{\theta}(\mathbf{z}_t|\mathbf{z}_{<t})$ can be defined as (e.g. with an LSTM)

$$p_{\theta}(\mathbf{z}_t|\mathbf{z}_{<t}) = p_{\theta}(\mathbf{z}_t|\mathbf{h}_t^p, \mathbf{c}_t^p), \quad \mathbf{h}_t^p, \mathbf{c}_t^p = \text{LSTM}_{\theta}(\mathbf{z}_{<t}). \quad (40)$$

This means the previous latent states $\mathbf{z}_{<t}$ are summarised by the LSTM internal recurrent states \mathbf{h}_t^p and \mathbf{c}_t^p , which are then transformed into the distributional parameters of $p_{\theta}(\mathbf{z}_t|\mathbf{z}_{<t})$. For the factorised encoder distribution, it can also be defined using an LSTM:

$$q_{\phi}(\mathbf{z}_t|\mathbf{x}_{\leq t}) = q_{\phi}(\mathbf{z}_t|\mathbf{h}_t^q, \mathbf{c}_t^q), \quad \mathbf{h}_t^q, \mathbf{c}_t^q = \text{LSTM}_{\phi}(\mathbf{x}_{\leq t}). \quad (41)$$

References

- Arjovsky, M., Shah, A., and Bengio, Y. (2016). Unitary evolution recurrent neural networks. In *International Conference on Machine Learning*, pages 1120–1128.
- Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A., Jozefowicz, R., and Bengio, S. (2016). Generating sentences from a continuous space. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 10–21. Association for Computational Linguistics.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734. Association for Computational Linguistics.
- Fabius, O. and van Amersfoort, J. R. (2014). Variational recurrent auto-encoders. *arXiv preprint arXiv:1412.6581*.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.

- Le, Q. V., Jaitly, N., and Hinton, G. E. (2015). A simple way to initialize recurrent networks of rectified linear units. *arXiv preprint arXiv:1504.00941*.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems*, 26:3111–3119.
- Pennington, J., Socher, R., and Manning, C. D. (2014). GloVe: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. (2014). Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *International Conference on Learning Representations*.
- Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27:3104–3112.