

Lecture 14. Recurrent Neural Networks and Continuous-depth Neural Networks

Bao Wang

Department of Mathematics

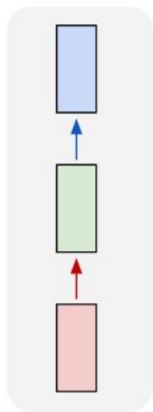
Scientific Computing and Imaging Institute

University of Utah

Math 5750/6880, Fall 2021

Feed-forward neural network: Recap

one to one



$$f(\mathbf{x}, \mathbf{W}) = \mathbf{W}_n \sigma \left(\mathbf{W}_{n-1} \sigma \left(\cdots \mathbf{W}_2 \sigma \left(\mathbf{W}_1 \mathbf{x} \right) \cdots \right) \right)$$

Image classification.

Recurrent Neural Networks

Learning sequences with varying sequence length!

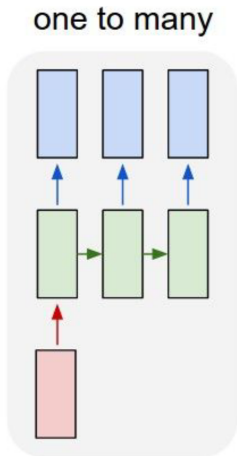
Why existing ConvNets are insufficient?

Variable sequence length inputs and outputs!

How to build a machine learning model for variable sequence length inputs and outputs?

Recurrent neural networks – I (Process sequences)

One to many



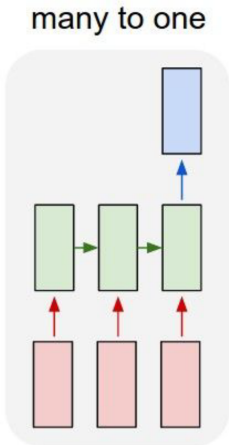
Man Is Walking Down Snowy Hill



Image captioning (image \rightarrow sequence of words)

Recurrent neural networks – II (Process sequences)

Many to one



Review (X)

"This movie is fantastic! I really like it because it is so good!"

"Not to my taste, will skip and watch another movie"

"This movie really sucks! Can I get my money back please?"

Rating (Y)

★★★★☆

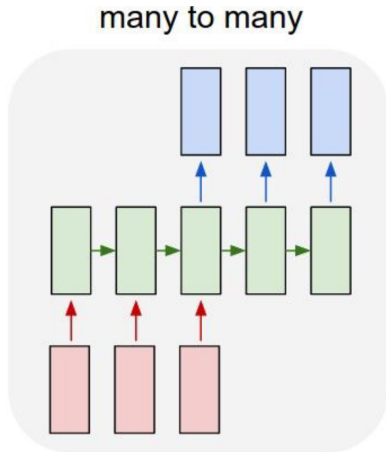
★★☆☆☆

★☆☆☆☆

Sentiment classification/sequence of video frames → action class

Recurrent neural networks – III (Process sequences)

Many to many

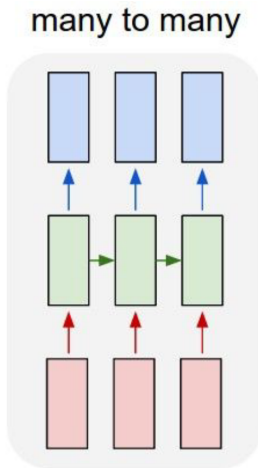


Machine translation, dialogue

Sequence of video frames → caption

Recurrent neural networks – IV (Process sequences)

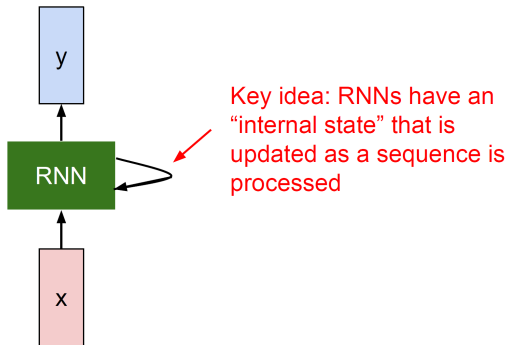
Many to many



Video captioning (sequence of video frames \rightarrow caption)

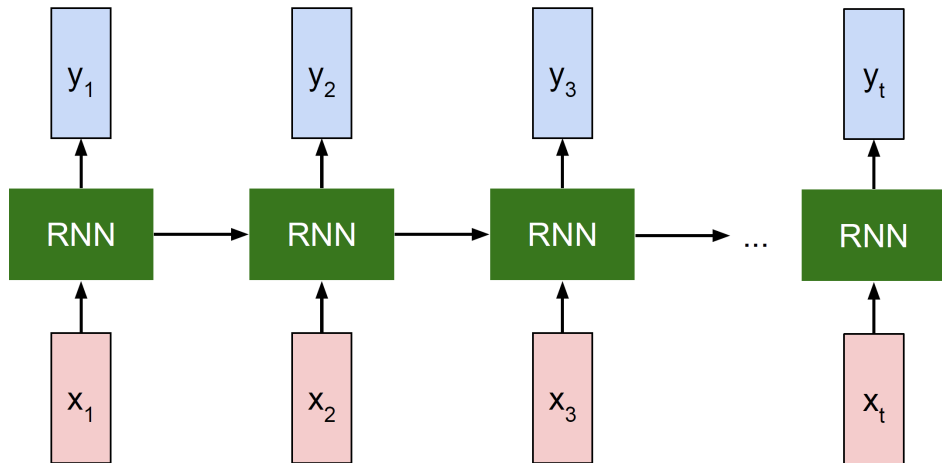
Recurrent neural networks

Let us start with a task that takes a variable input and produces an output at every step. **Many-to-many.**



RNN internal state: ***h***.

Unrolled RNN



RNN internal state: h_t is a function of h_{t-1} and x_t .

RNN hidden state update

We can process a sequence of vectors $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots]$ by applying a recurrence formula at every time step:

$$\mathbf{h}_t = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t),$$

where

- \mathbf{h}_t is the new state
- $f_{\mathbf{W}}(\cdot, \cdot)$ is a function with parameters \mathbf{W}
- \mathbf{h}_{t-1} is the old state
- \mathbf{x}_t is the input vector at some time step

RNN output generation

We output \mathbf{y}_t at every time-step t by using the following formula

$$\mathbf{y}_t = f_{\mathbf{W}_o}(\mathbf{h}_t),$$

where

- \mathbf{y}_t is the output
- $f_{\mathbf{W}_o}(\cdot)$ is a function with parameters \mathbf{W}_o
- \mathbf{h}_t is the new state

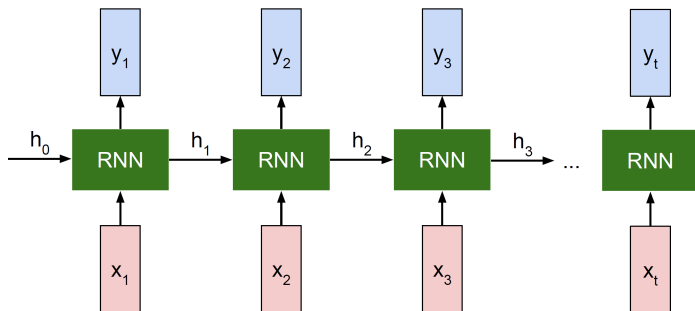
RNN

We can process a sequence of vectors $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots]$ by applying a recurrence formula at every time step:

$$\mathbf{h}_t = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t),$$

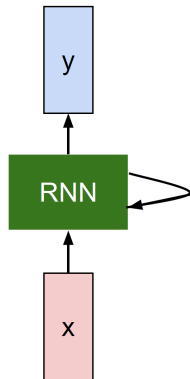
output \mathbf{y}_t at time-step t by using

$$\mathbf{y}_t = f_{\mathbf{W}_o}(\mathbf{h}_t).$$



RNNs

$$\mathbf{h}_t = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t)$$



Note that the same function and the same set of parameters are used at every time step.

Simple recurrent neural network

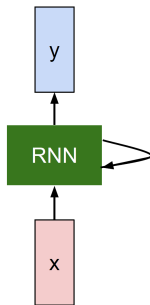
Let $\mathbf{h}_t = f_{\mathbf{W}}(\mathbf{h}_{t-1}, \mathbf{x}_t)$ be

$$\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t).$$

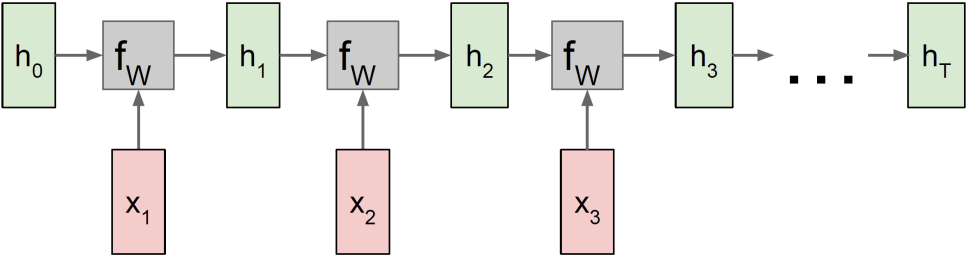
Therefore the simple RNN model can be formulated as

$$\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t)$$

$$\mathbf{y}_t = \mathbf{W}_{hy}\mathbf{h}_t \text{ or } \mathbf{y}_t = \text{softmax}(\mathbf{W}_{hy}\mathbf{h}_t)$$

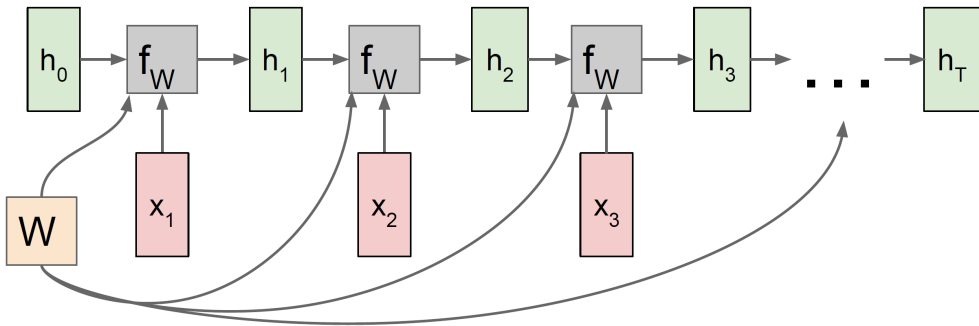


RNN: Computational graph



$$h_t = f_W(h_{t-1}, x_t)$$

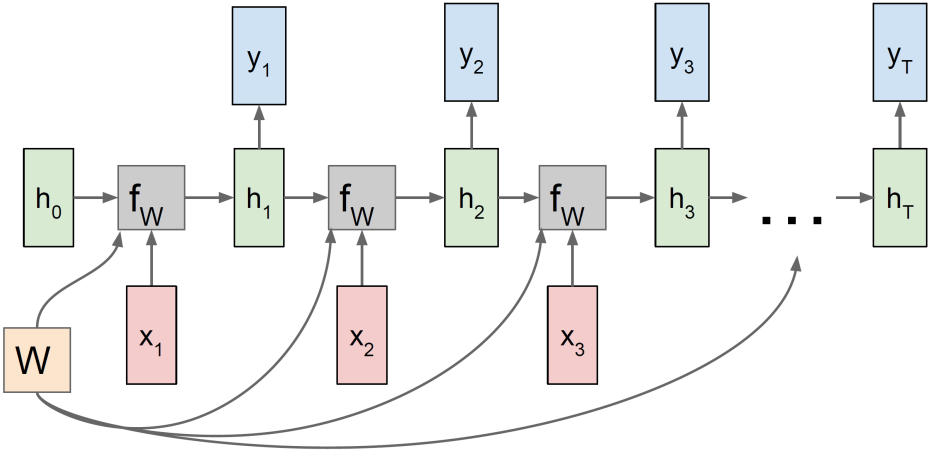
RNN: Computational graph



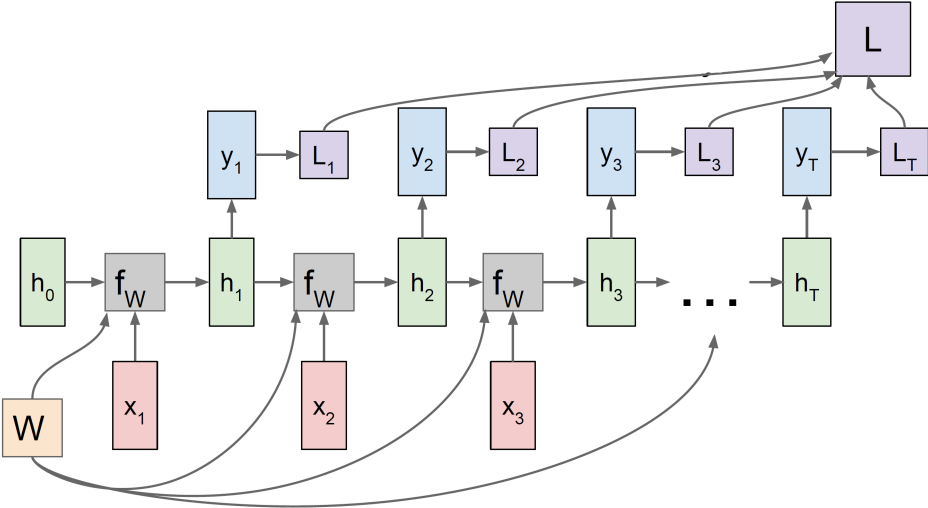
$$h_t = f_W(h_{t-1}, x_t)$$

Re-use the same weight matrix at every time-step.

RNN: Computational graph: Many-to-many

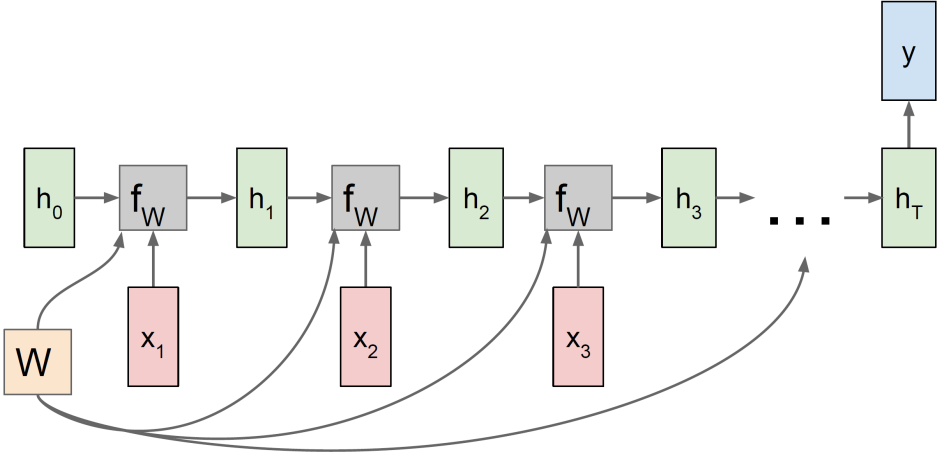


RNN: Computational graph: Many-to-many

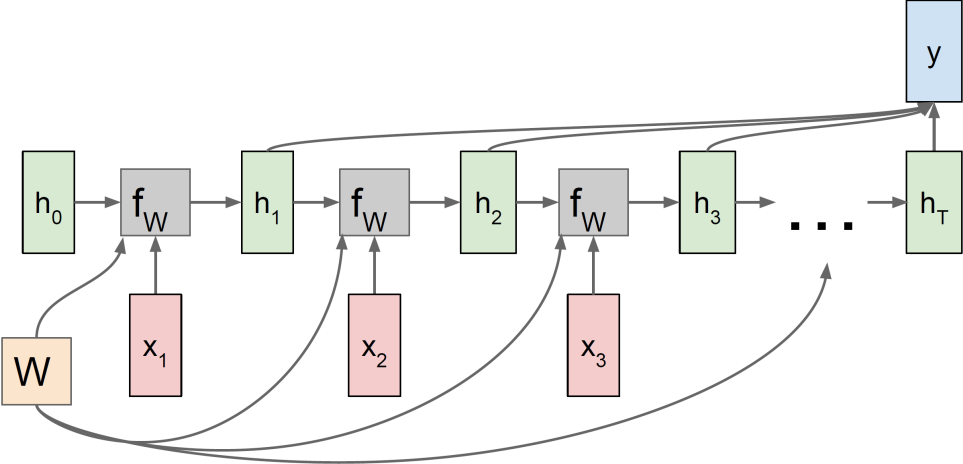


The loss is the sum of the loss at each slot.

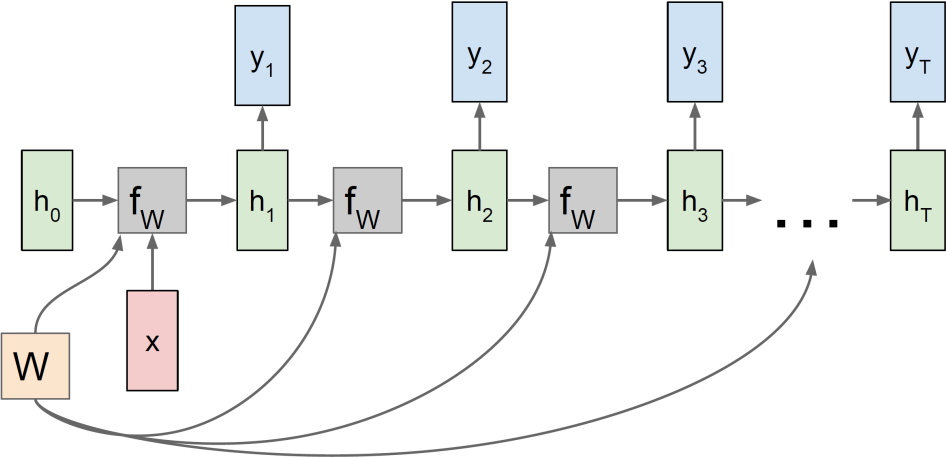
RNN: Computational graph: Many-to-one – I



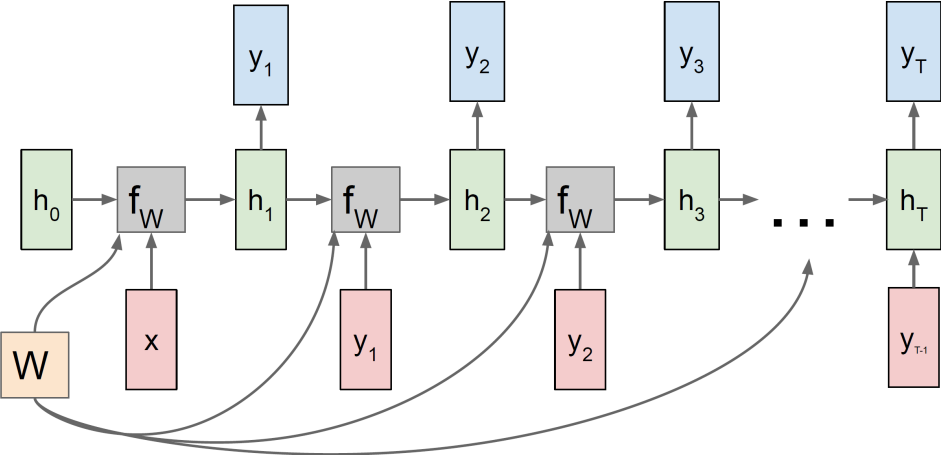
RNN: Computational graph: Many-to-one – II



RNN: Computational graph: One-to-many – I

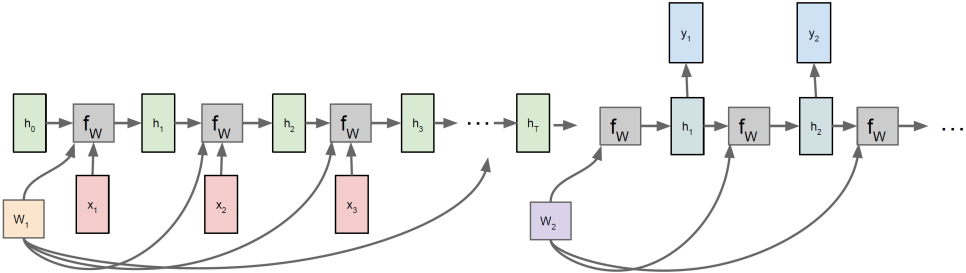


RNN: Computational graph: One-to-many – II



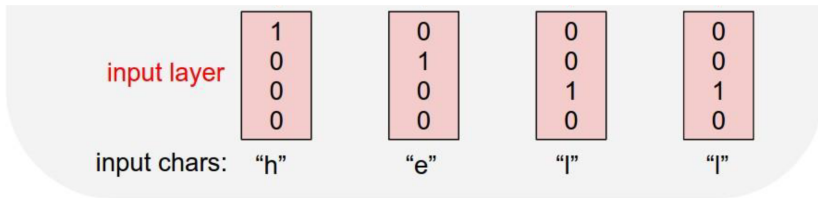
Sequence to sequence modeling: Many-to-one + One-to-many

- **Many-to-one:** Encode input sequence (x_1, \dots, x_T) in a single vector h_T .
- **One-to-many:** Produce output sequence (y_1, y_2, \dots) from single input vector h_T .

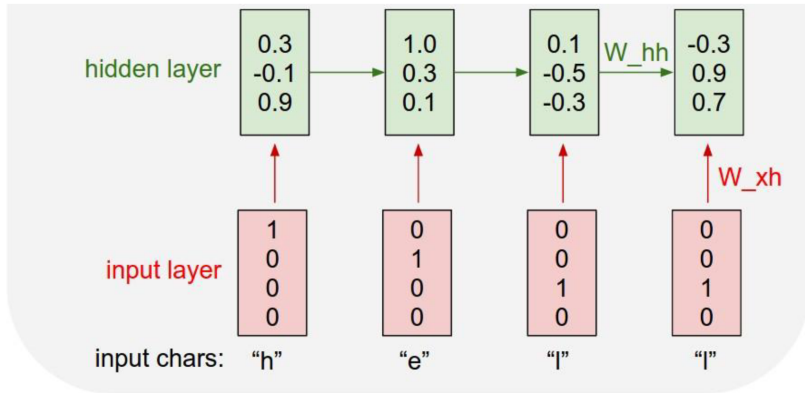


Example: Character-level language model

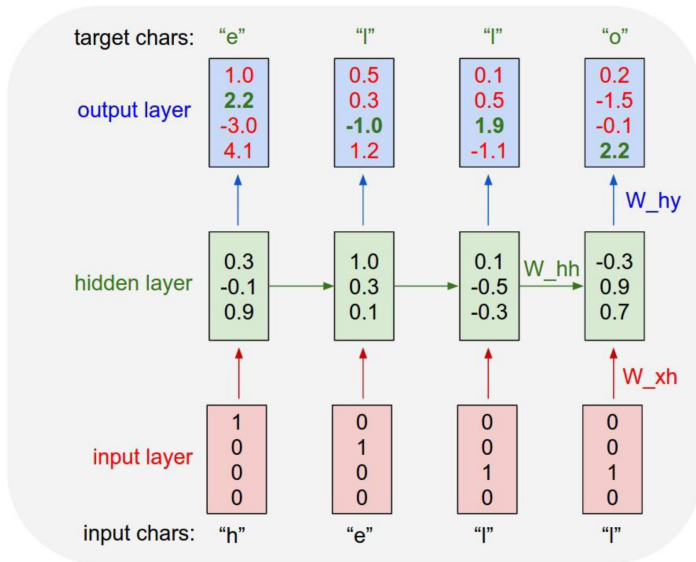
- Vocabulary: [h,e,l,o]
- Example training sequence: "hello"
- Task: predict the next character



$$\mathbf{h}_t = \tanh \left(\mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{W}_{xh} \mathbf{x}_t \right)$$

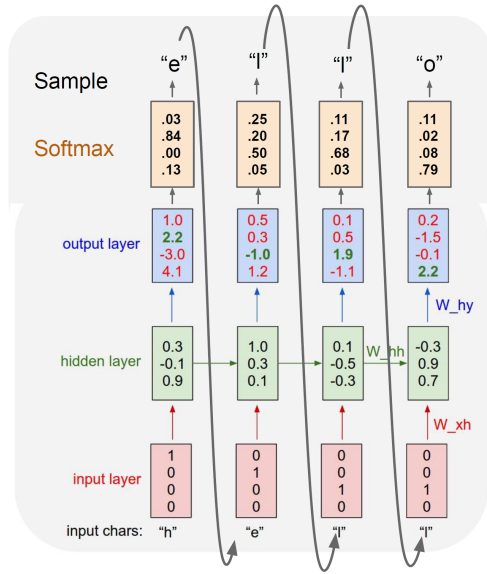


Training



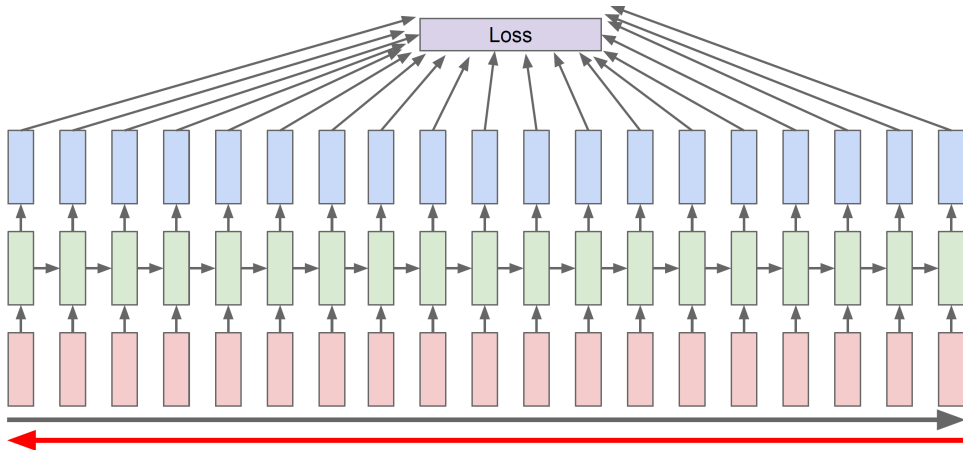
Testing

At test-time sample characters one at a time, and then feed back to model to predict the next.



Backpropagation through time

Forward through entire sequence to compute loss, then backward through entire sequence to compute gradient. (We will discuss the algorithm later).



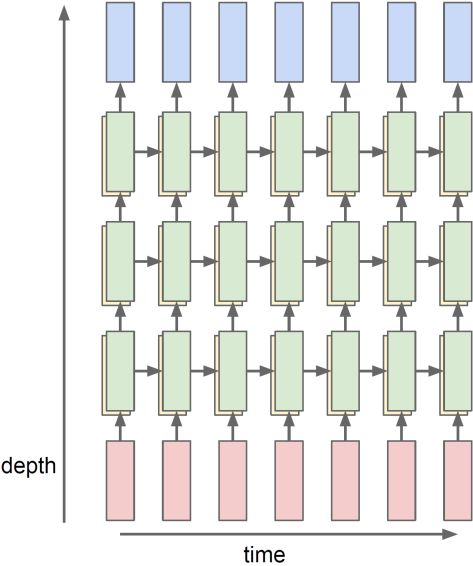
RNN Advantages:

- Can process any length input
- Computation for step t can (in theory) use information from many steps back
- Model size does not increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

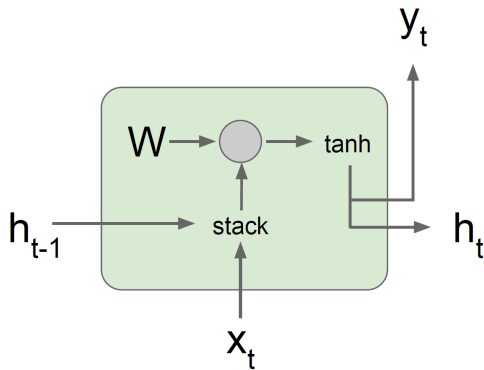
RNN Disadvantages:

- Recurrent computation is slow
- In practice, difficult to access information from many steps back

Multilayer RNNs

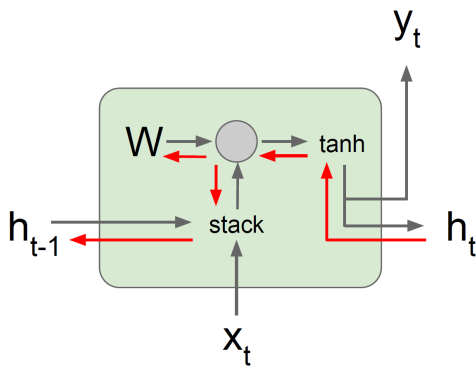


$$\mathbf{h}_t = \tanh \left(\mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{W}_{hx} \mathbf{x}_t \right) = \tanh \left(\begin{pmatrix} \mathbf{W}_{hh} & \mathbf{W}_{hx} \end{pmatrix} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{pmatrix} \right) := \tanh \left(\mathbf{W} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{pmatrix} \right)$$



BPTT

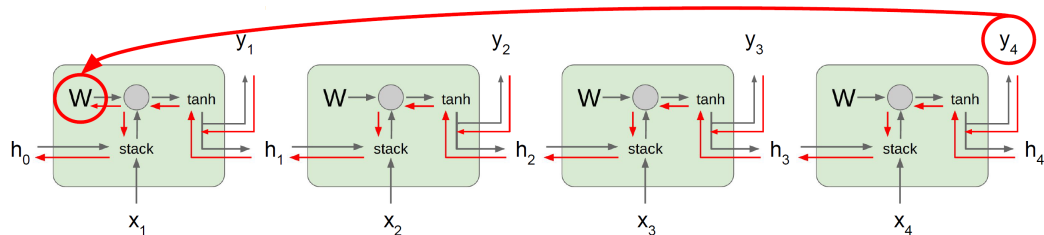
$$\mathbf{h}_t = \tanh(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t) = \tanh\left(\begin{pmatrix} \mathbf{W}_{hh} & \mathbf{W}_{hx} \end{pmatrix} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{pmatrix}\right) := \tanh\left(\mathbf{W} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{pmatrix}\right)$$



$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \tanh'(\mathbf{W}_{hh}\mathbf{h}_{t-1} + \mathbf{W}_{xh}\mathbf{x}_t) \mathbf{W}_{hh}$$

Backpropagation from \mathbf{h}_t to \mathbf{h}_{t-1} multiplies by \mathbf{W} (actually \mathbf{W}_{hh}^\top).

BPTT



Note that

$$\frac{\partial L}{\partial \mathbf{W}} = \sum_{t=1}^T \frac{\partial L_t}{\partial \mathbf{W}},$$

where

$$\frac{\partial L_T}{\partial \mathbf{W}} = \frac{\partial L_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \cdots \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} = \frac{\partial L_T}{\partial \mathbf{h}_T} \left(\prod_{t=2}^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right) \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}}$$

Vanishing and exploding gradient issue in BPTT

$$\frac{\partial L_T}{\partial \mathbf{W}} = \frac{\partial L_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \cdots \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} = \frac{\partial L_T}{\partial \mathbf{h}_T} \left(\prod_{t=2}^T \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} \right) \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}}.$$

Since

$$\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \tanh' \left(\mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{W}_{xh} \mathbf{x}_t \right) \mathbf{W}_{hh},$$

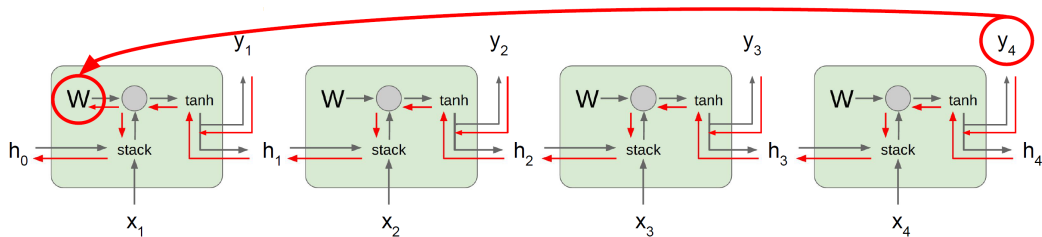
therefore

$$\frac{\partial L_T}{\partial \mathbf{W}} = \frac{\partial L_T}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_{T-1}} \cdots \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}} = \frac{\partial L_T}{\partial \mathbf{h}_T} \left(\prod_{t=2}^T \tanh' \left(\mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{W}_{xh} \mathbf{x}_t \right) \right) \mathbf{W}_{hh}^{T-1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}}.$$

Note that $\tanh' \left(\mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{W}_{xh} \mathbf{x}_t \right) < 1$ a.s. Hence, $\frac{\partial L_T}{\partial \mathbf{W}} \rightarrow 0$ as $T \rightarrow \infty$.

Vanishing gradient: bottleneck for learning long-term dependencies.

Vanishing and exploding gradient issue in BPTT



If we assume there is no nonlinearity, we have

$$\frac{\partial L_T}{\partial \mathbf{W}} = \frac{\partial L_T}{\partial \mathbf{h}_T} \mathbf{W}_{hh}^{T-1} \frac{\partial \mathbf{h}_1}{\partial \mathbf{W}}.$$

- Largest singular value $> 1 \Rightarrow$ exploding gradients. [Gradient clipping, i.e., scaling gradient if its norm is greater than 1.]
- Largest singular value $< 1 \Rightarrow$ vanishing gradients. [Change the RNN architecture.]

Long-short Term Memory (LSTM)

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997.

LSTM

Vanilla RNN:

$$\mathbf{h}_t = \tanh \left(\mathbf{W} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{pmatrix} \right)$$

LSTM:

$$\begin{pmatrix} \mathbf{i} \\ \mathbf{f} \\ \mathbf{o} \\ \mathbf{g} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \mathbf{W} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{pmatrix},$$

$$\mathbf{c}_t = \mathbf{f} \odot \mathbf{c}_{t-1} + \mathbf{i} \odot \mathbf{g}$$

$$\mathbf{h}_t = \mathbf{o} \odot \tanh(\mathbf{c}_t)$$

where \odot represents entry-wise product.

LSTM

$$\begin{pmatrix} \mathbf{i} \\ \mathbf{f} \\ \mathbf{o} \\ \mathbf{g} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \mathbf{W} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{pmatrix},$$

$$\mathbf{c}_t = \mathbf{f} \odot \mathbf{c}_{t-1} + \mathbf{i} \odot \mathbf{g}$$

$$\mathbf{h}_t = \mathbf{o} \odot \tanh(\mathbf{c}_t)$$

where

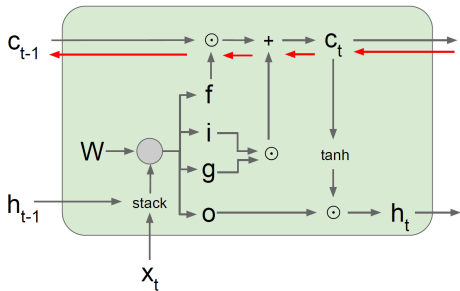
- \mathbf{i} : input gate, whether to write to cell
- \mathbf{f} : forget gate, whether to erase cell
- \mathbf{o} : output gate, how much to reveal cell
- \mathbf{g} How much to write to cell

LSTM

$$\begin{pmatrix} \mathbf{i} \\ \mathbf{f} \\ \mathbf{o} \\ \mathbf{g} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \mathbf{W} \begin{pmatrix} \mathbf{h}_{t-1} \\ \mathbf{x}_t \end{pmatrix},$$

$$\mathbf{c}_t = \mathbf{f} \odot \mathbf{c}_{t-1} + \mathbf{i} \odot \mathbf{g}$$

$$\mathbf{h}_t = \mathbf{o} \odot \tanh(\mathbf{c}_t)$$



Remarks

LSTM does not guarantee that there is no vanishing/exploding gradient, but it does provide an easier way for the model to learn long-term dependencies in practice.

Momentum Recurrent Neural Networks

T. Nguyen, R. Baraniuk, A. Bertozzi, S. Osher, B. Wang, MomentumRNN: Integrating momentum into RNNs, NeurIPS, 2020.

Heavy-ball method vs. RNN

Heavy-ball method

$$\mathbf{x}_t = \mathbf{x}_{t-1} - s\nabla f(\mathbf{x}_{t-1}) + \mu(\mathbf{x}_{t-1} - \mathbf{x}_{t-2}),$$

let $\mathbf{p}_t = \mathbf{x}_{t-1} - \mathbf{x}_t$, we can rewrite the heavy-ball method as

$$\mathbf{p}_t = \mu\mathbf{p}_{t-1} + s\nabla f(\mathbf{x}_{t-1}); \quad \mathbf{x}_t = \mathbf{x}_{t-1} - \mathbf{p}_t.$$

Heavy-ball method vs. RNN

Heavy-ball method:

$$\mathbf{p}_t = \mu \mathbf{p}_{t-1} + s \nabla f(\mathbf{x}_{t-1}); \quad \mathbf{x}_t = \mathbf{x}_{t-1} - \mathbf{p}_t.$$

Also, the recurrence equation of RNN can be written as

$$\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{x}_t),$$

let $\phi(\cdot) = \sigma(\mathbf{U}(\cdot))$ and $\mathbf{u}_t = \mathbf{U}^{-1}\mathbf{W}\mathbf{x}_t$, we can rewrite the above equation as

$$\mathbf{h}_t = \phi(\mathbf{h}_{t-1} + \mathbf{u}_t).$$

What is we treat $-\mathbf{u}_t$ in RNN as $\nabla f(\mathbf{x}_t)$ in the heavy-ball method? Also, regard $\phi(\cdot)$ as a projection.

Momentum RNN

$$\mathbf{p}_t = \mu \mathbf{p}_{t-1} - s \mathbf{u}_t; \quad \mathbf{h}_t = \phi(\mathbf{h}_{t-1} - \mathbf{p}_t),$$

where $\mu \geq 0, s > 0$ are two hyperparameters.

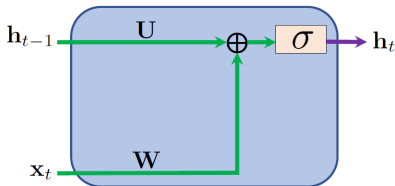
Let $\mathbf{v}_t := -\mathbf{U}\mathbf{p}_t$, we arrive at the following Momentum RNN

$$\mathbf{v}_t = \mu \mathbf{v}_{t-1} + s \mathbf{W}\mathbf{x}_t; \quad \mathbf{h}_t = \sigma(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{v}_t).$$

Momentum RNN

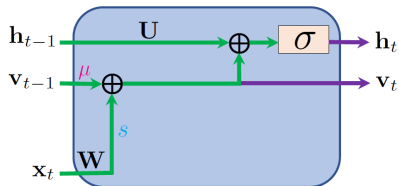
Recurrent Cell

$$\mathbf{h}_t = \phi(\mathbf{U} \times \mathbf{h}_{t-1} + \underbrace{\mathbf{W} \times \mathbf{x}_t}_{\text{gradient}})$$



Momentum Cell

$$\mathbf{v}_t = \underbrace{\mu}_{\text{momentum}} \times \mathbf{v}_{t-1} + \underbrace{s}_{\text{stepsize}} \times \mathbf{W} \times \mathbf{x}_t$$
$$\mathbf{h}_t = \sigma(\mathbf{U} \times \mathbf{h}_{t-1} + \mathbf{v}_t)$$



Momentum RNN overcomes vanishing gradients

RNN BPTT:

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \prod_{k=t}^{T-1} \frac{\partial \mathbf{h}_{k+1}}{\partial \mathbf{h}_k} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \prod_{k=t}^{T-1} (\mathbf{D}_k \mathbf{U}^\top),$$

where $\mathbf{D}_k = \text{diag}(\sigma'(\mathbf{U}\mathbf{h}_k + \mathbf{W}\mathbf{x}_{k+1}))$ is a diagonal matrix with $\sigma'(\mathbf{U}\mathbf{h}_k + \mathbf{W}\mathbf{x}_{k+1})$ being its diagonal entries.

Momentum RNN overcomes vanishing gradients

Momentum RNN:

$$\mathbf{v}_t = \mu \mathbf{v}_{t-1} + s \mathbf{W} \mathbf{x}_t; \quad \mathbf{h}_t = \sigma(\mathbf{U} \mathbf{h}_{t-1} + \mathbf{v}_t),$$

which can be rewritten as

$$\mathbf{h}_t = \sigma \left(\mathbf{U}(\mathbf{h}_{t-1} - \mu \mathbf{h}_{t-2}) + \mu \sigma^{-1}(\mathbf{h}_{t-1}) + s \mathbf{W} \mathbf{x}_t \right),$$

where $\sigma^{-1}(\cdot)$ is the inverse function of $\sigma(\cdot)$. We compute $\partial L / \partial \mathbf{h}_t$ as follows

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \prod_{k=t}^{T-1} \frac{\partial \mathbf{h}_{k+1}}{\partial \mathbf{h}_k} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \prod_{k=t}^{T-1} \hat{\mathbf{D}}_k [\mathbf{U}^\top + \mu \mathbf{\Sigma}_k], \quad (1)$$

where $\hat{\mathbf{D}}_k = \text{diag}(\sigma'(\mathbf{U}(\mathbf{h}_k - \mu \mathbf{h}_{k-1}) + \mu \sigma^{-1}(\mathbf{h}_k) + s \mathbf{W} \mathbf{x}_{k+1}))$ and $\mathbf{\Sigma} = \text{diag}((\sigma^{-1})'(\mathbf{h}_k))$. For mostly used σ , e.g., sigmoid and tanh, $(\sigma^{-1}(\cdot))' > 1$ and $\mu \mathbf{\Sigma}_k$ dominates \mathbf{U}^\top .

Neural Ordinary Differential Equations

R. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, Neural Ordinary Differential Equations, NeurIPS 2018.

Neural ODEs

Neural ODE:

$$\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta),$$

where $f(\mathbf{h}(t), t, \theta)$ is a neural network, e.g., $f(\mathbf{h}(t), t, \theta) = \theta_2 \sigma(\theta_1 \mathbf{h}(t))$.

Forward propagation: starting from the input $\mathbf{h}(0)$, we define the output layer $\mathbf{h}(T)$ to be the solution to the above ODE initial value problem at some time T . $\mathbf{h}(T)$ can be computed by a black-box numerical ODE solver, which evaluates the hidden unit dynamics f whenever necessary to determine the solution with the desired accuracy.

Loss function: consider the scalar-valued loss function $L(\cdot)$, whose input is the result of an ODE solver

$$L(\mathbf{h}(T)) = L\left(h(0) + \int_0^T f(\mathbf{h}(t), t, \theta) dt\right) = L(\text{ODESolve}(\mathbf{h}(0), f, t_0, t_1, \theta)).$$

Numerical ODE solver – Forward Euler method

$$\mathbf{h}_{k+1} = \mathbf{h}_k + sf(\mathbf{h}_k, t_k, \theta)$$

Training neural ODEs: Adjoint sensitivity method

Adjoint state: $\mathbf{a}(t) := \frac{\partial L}{\partial \mathbf{h}(t)}$

Theorem 1.

$$\frac{dL}{d\theta} = - \int_T^0 \mathbf{a}(t)^\top \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \theta} dt,$$

where $\mathbf{a}(t)$ satisfies the following adjoint ODE

$$\frac{d\mathbf{a}(t)}{dt} = -\mathbf{a}(t)^\top \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}}.$$

Proof. Let $\mathbf{h}(t)$ follow the differential equation $\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta)$, where θ are the parameters. We will prove that if we define an adjoint state

$$\mathbf{a}(t) := \frac{dL}{d\mathbf{h}(t)}, \quad (2)$$

then it follows the following ODE

$$\frac{d\mathbf{a}(t)}{dt} = -\mathbf{a}(t) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}(t)}. \quad (3)$$

For ease of notation, we denote vectors as row vectors, whereas the vector in the previous part are column vectors.

The adjoint state is the gradient with respect to the hidden state at a specified time t . In standard neural networks, the gradient of a hidden layer \mathbf{h}_t depends on the gradient from the next layer \mathbf{h}_{t+1} by chain rule

$$\frac{dL}{d\mathbf{h}_t} = \frac{dL}{d\mathbf{h}_{t+1}} \frac{d\mathbf{h}_{t+1}}{d\mathbf{h}_t}. \quad (4)$$

With a continuous hidden state, we can write the transformation after an ϵ change in time as (note $\dot{\mathbf{h}}(t) = f(\mathbf{h}(t), t, \theta)$, implying $\mathbf{h}(t) = \int f(\mathbf{h}(t), t, \theta) dt$)

$$\mathbf{h}(t + \epsilon) = \int_t^{t+\epsilon} f(\mathbf{h}(t), t, \theta) dt + \mathbf{h}(t) := T_\epsilon(\mathbf{h}(t), t) \quad (5)$$

and chain rule can also be applied ($\partial\mathbf{h}(t + \epsilon)/\partial\mathbf{h}(t) = \partial T_\epsilon(\mathbf{h}(t), t)/\partial\mathbf{h}(t)$)

$$\frac{dL}{d\mathbf{h}(t)} = \frac{dL}{d\mathbf{h}(t + \epsilon)} \frac{d\mathbf{h}(t + \epsilon)}{d\mathbf{h}(t)} \quad \text{or} \quad \mathbf{a}(t) = \mathbf{a}(t + \epsilon) \frac{\partial T_\epsilon(\mathbf{h}(t), t)}{\partial\mathbf{h}(t)}. \quad (6)$$

The proof of (3) follows from the definition of derivative:

$$\begin{aligned}
 & \frac{d\mathbf{a}(t)}{dt} \\
 &= \lim_{\epsilon \rightarrow 0^+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t)}{\epsilon} \\
 &= \lim_{\epsilon \rightarrow 0^+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t + \epsilon) \partial(T_\epsilon(\mathbf{h}(t))) / (\partial \mathbf{h}(t))}{\epsilon} \quad (\text{by (6)}) \\
 &= \lim_{\epsilon \rightarrow 0^+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t + \epsilon) \frac{\partial}{\partial \mathbf{h}(t)} (\mathbf{h}(t) + \epsilon f(\mathbf{h}(t), t, \theta) + \mathcal{O}(\epsilon^2))}{\epsilon} \quad (\text{Taylor series of } \mathbf{h}(t + \epsilon) = T_\epsilon(\mathbf{h}(t), \epsilon)) \\
 &= \lim_{\epsilon \rightarrow 0^+} \frac{\mathbf{a}(t + \epsilon) - \mathbf{a}(t + \epsilon) \left(\mathbf{I} + \epsilon \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}(t)} + \mathcal{O}(\epsilon^2) \right)}{\epsilon} \tag{7} \\
 &= \lim_{\epsilon \rightarrow 0^+} \frac{-\epsilon \mathbf{a}(t + \epsilon) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}(t)} + \mathcal{O}(\epsilon^2)}{\epsilon} \\
 &= \lim_{\epsilon \rightarrow 0^+} -\mathbf{a}(t + \epsilon) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}(t)} + \mathcal{O}(\epsilon) \\
 &= -\mathbf{a}(t) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}(t)}
 \end{aligned}$$

Next, we compute gradient w.r.t. θ . We can generalize (3) to obtain gradient w.r.t. θ and t . We view θ ($\theta(t) = \theta$) and t ($t(t) = t$) as states with constant differential equations and write

$$\frac{\partial \theta(t)}{\partial t} = 0, \quad \frac{dt(t)}{dt} = 1. \quad (8)$$

We can then combine these with \mathbf{h} to form an augmented state with corresponding differential equation and adjoint state,

$$\frac{d}{dt} \begin{pmatrix} \mathbf{h} \\ \theta \\ t \end{pmatrix} (t) = f_{aug}([\mathbf{h}, \theta, t]) := \begin{pmatrix} f([\mathbf{h}, \theta, t]) \\ 0 \\ 1 \end{pmatrix}, \quad (9)$$

where

$$\mathbf{a}_{aug} := \begin{pmatrix} \mathbf{a} \\ \mathbf{a}_\theta \\ \mathbf{a}_t \end{pmatrix}, \quad \mathbf{a}_\theta(t) := \frac{dL}{d\theta(t)}, \quad \mathbf{a}_t(t) := \frac{dL}{dt(t)}.$$

The Jacobian of f has the form

$$\frac{\partial f_{aug}(t)}{\partial[\mathbf{h}, \theta, t]} = \begin{pmatrix} \frac{\partial f}{\partial \mathbf{h}} & \frac{\partial f}{\partial \theta} & \frac{\partial f}{\partial t} \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} (t) \quad (10)$$

where 0 is a matrix of zeros with the appropriate dimensions. We plug this into (3) to obtain

$$\frac{d\mathbf{a}_{aug}(t)}{dt} = -[\mathbf{a}(t) \ \mathbf{a}_\theta(t) \ \mathbf{a}_t(t)] \frac{\partial f_{aug}}{\partial[\mathbf{h}, \theta, t]}(t) = - \left[\mathbf{a} \frac{\partial f}{\partial \mathbf{h}} \ \mathbf{a} \frac{\partial f}{\partial \theta} \ \mathbf{a} \frac{\partial f}{\partial t} \right] (t).$$

The first element is the adjoint differential equation (3), as expected. The second element can be used to obtain the total gradient w.r.t. the parameters, by integrating over the full interval and setting $\mathbf{a}_\theta(T) = 0$ (L is independent of $\theta(T)$).

$$\frac{dL}{d\theta} = \mathbf{a}_\theta(0) = - \int_T^0 \mathbf{a}(t) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \theta} dt. \quad (11)$$

Advantages of Neural ODEs

- **Memory efficiency:** Recall that in backpropagation we need to store any intermediate quantities of the forward pass, which is expensive in memory footprint. Later, we will see that in training neural ODEs we do not need to store these intermediate quantities.
- **Continuous time-series models:** Neural ODE can learn time-series that are irregularly-observed in time, which is significantly different from RNNs.

Computational bottleneck of neural ODEs

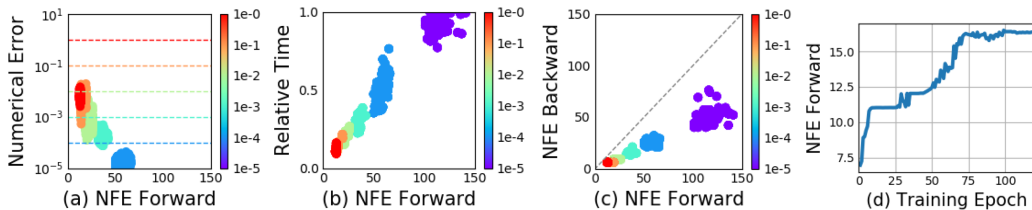


Figure 3: Statistics of a trained ODE-Net. (NFE = number of function evaluations.)

NFEs can be excessively high in training neural ODEs, resulting in complicated models.

As training goes, the stiffness of NODE keeps increasing!

Learning bottleneck of neural ODEs

Lemma 1. For neural ODEs, we have

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \exp \left\{ - \int_T^t \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}(s), s, \theta) ds \right\},$$

the adjoint state $\partial L / \partial \mathbf{h}_t$ vanishes quickly, making neural ODE cannot learn long-term dependencies.

Note that if $\int_T^t \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}(s), s, \theta)$ is large, then $\frac{\partial L}{\partial \mathbf{h}_t} \approx 0$, i.e., vanishing gradient.

Remark. BPTT for training the RNN model $\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{x}_t)$:

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \prod_{k=t}^{T-1} \frac{\partial \mathbf{h}_{k+1}}{\partial \mathbf{h}_k} = \frac{\partial L}{\partial \mathbf{h}_T} \cdot \prod_{k=t}^{T-1} (\mathbf{D}_k \mathbf{U}^\top),$$

where $\mathbf{D}_k = \text{diag}(\sigma'(\mathbf{U}\mathbf{h}_k + \mathbf{W}\mathbf{x}_{k+1}))$ is a diagonal matrix with $\sigma'(\mathbf{U}\mathbf{h}_k + \mathbf{W}\mathbf{x}_{k+1})$ being its diagonal entries.

Heavy-ball Neural Ordinary Differential Equations

H. Xia*, V. Suliafu*, H. Ji, T. Nguyen, A. Bertozzi, S. Osher, B. Wang, Heavy-ball Neural Ordinary Differential Equations, NeurIPS 2021.

Recap: Derivation of the heavy-ball ODE

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) + \theta(\mathbf{x}^t - \mathbf{x}^{t-1}),$$

Let $\mathbf{m}^{t+1} := (\mathbf{x}^{t+1} - \mathbf{x}^t)/\sqrt{\eta}$ and let $\theta := 1 - \gamma\sqrt{\eta}$, where $\gamma \geq 0$ is another hyperparameter. Then we have

$$\frac{\mathbf{x}^{t+1} - \mathbf{x}^t}{\sqrt{\eta}} = -\sqrt{\eta} \nabla f(\mathbf{x}^t) + (1 - \gamma\sqrt{\eta}) \frac{\mathbf{x}^t - \mathbf{x}^{t-1}}{\sqrt{\eta}},$$

therefore, we can rewrite the heavy-ball method as

$$\mathbf{m}^{t+1} = (1 - \gamma\sqrt{\eta})\mathbf{m}^t - \sqrt{\eta} \nabla f(\mathbf{x}^t); \quad \mathbf{x}^{t+1} = \mathbf{x}^t + \sqrt{\eta} \mathbf{m}^{t+1}.$$

Let $\eta \rightarrow 0$; we obtain the following system of first-order ODEs

$$\frac{d\mathbf{X}(t)}{dt} = \mathbf{M}(t); \quad \frac{d\mathbf{M}(t)}{dt} = -\gamma\mathbf{M}(t) - \nabla f(\mathbf{X}(t)),$$

which can be further written as

$$\ddot{\mathbf{X}}(\tau) + \gamma\dot{\mathbf{X}}(\tau) + \nabla f(\mathbf{X}(\tau)) = 0.$$

What if we parameterize $-\nabla f(\mathbf{X}(\tau))$ by a neural network?

Heavy-ball neural ODE can accelerate forward and backward propagation.

Heavy-ball Neural ODEs

Heavy-ball neural ODE (HBNODE)

$$\frac{d^2 \mathbf{h}(t)}{dt^2} + \gamma \frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta), \quad \gamma \geq 0. \quad (12)$$

The above HBNODE is equivalent to the following system of first-order neural ODEs

$$\frac{d\mathbf{h}(t)}{dt} = \mathbf{m}(t); \quad \frac{d\mathbf{m}(t)}{dt} = -\gamma \mathbf{m}(t) + f(\mathbf{h}(t), t, \theta). \quad (13)$$

Stiffness analysis of the linearized models

Let us compare the stiffness of the following two ODE models

$$\frac{d\mathbf{h}(t)}{dt} = \mathbf{A}\mathbf{h}(t)$$

and

$$\begin{aligned} \frac{d\mathbf{h}(t)}{dt} &= \mathbf{m}(t) \\ \frac{d\mathbf{m}(t)}{dt} &= -\gamma\mathbf{m}(t) + \mathbf{A}\mathbf{h}(t), \end{aligned} \quad \text{i.e.,} \quad \frac{d}{dt} \begin{pmatrix} \mathbf{h}(t) \\ \mathbf{m}(t) \end{pmatrix} = \underbrace{\begin{pmatrix} 0 & \mathbf{I} \\ \mathbf{A} & -\gamma\mathbf{I} \end{pmatrix}}_{:=\mathbf{B}} \begin{pmatrix} \mathbf{h}(t) \\ \mathbf{m}(t) \end{pmatrix}$$

The stiffness of ODEs is defined by the ratio between the largest and the smallest eigenvalues (in absolute value). We consider the case when \mathbf{A} is positive definite.

Stiffness of the linearized models

Let the eigenvalues and the corresponding eigenvectors of $\mathbf{A} \in \mathbb{R}^{d \times d}$ be $\{\lambda_i, \mathbf{v}_i\}_{i=1}^d$. We assume the eigenvectors of \mathbf{B} has the following form

$$\begin{pmatrix} \mathbf{v}_i \\ c\mathbf{v}_i \end{pmatrix}, \text{ where } c \text{ is a constant,}$$

then we have

$$\begin{pmatrix} 0 & \mathbf{I} \\ \mathbf{A} & -\gamma\mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{v}_i \\ c\mathbf{v}_i \end{pmatrix} = \begin{pmatrix} c\mathbf{v}_i \\ \lambda_i\mathbf{v}_i - \gamma c\mathbf{v}_i \end{pmatrix}$$

Let $\tilde{\lambda}$ be the eigenvalue of \mathbf{B} , then we have

$$\frac{c}{1} = \tilde{\lambda} = \frac{\lambda_i - \gamma c}{c}$$

Thus:

$$c^2 + \gamma c - \lambda_i = 0 \implies c = \tilde{\lambda} = \frac{\gamma \pm \sqrt{\gamma^2 + 4\lambda_i}}{2}$$

Stiffness of linearized models

Each eigenvalue λ_i of the matrix \mathbf{A} correspond to two eigenvalues of the matrix \mathbf{B} , given by

$$\tilde{\lambda}_{i,1} = \frac{\gamma + \sqrt{\gamma^2 + 4\lambda_i}}{2}, \quad \text{and} \quad \tilde{\lambda}_{i,2} = \frac{\gamma - \sqrt{\gamma^2 + 4\lambda_i}}{2}.$$

Stiffness of linearized models

We assume the eigenvalues of \mathbf{A} are sorted as follows:

$$0 < \lambda_{\min} \leq \lambda_2 \leq \lambda_3 \leq \dots \leq \lambda_{\max}.$$

Therefore, we have

$$\tilde{\lambda}_{\max} = \frac{\gamma + \sqrt{\gamma^2 + 4\lambda_{\max}}}{2} \quad \text{and} \quad |\tilde{\lambda}_{\min}| = \left| \frac{\gamma - \sqrt{\gamma^2 + 4\lambda_{\min}}}{2} \right|.$$

Thus the stiffness ratio of the matrix \mathbf{B} satisfies

$$S(\mathbf{B}) = \frac{\gamma + \sqrt{\gamma^2 + 4\lambda_{\max}}}{|\gamma - \sqrt{\gamma^2 + 4\lambda_{\min}}|} \underbrace{\leq}_{\text{let } \gamma=0} \sqrt{\frac{\lambda_{\max}}{\lambda_{\min}}} = \sqrt{S(\mathbf{A})}.$$

Analysis of the adjoint state of the heavy-ball neural ODEs.

Theorem 2. The adjoint state \mathbf{a}_h and \mathbf{a}_m of the following HBNODE

$$\frac{d\mathbf{h}(t)}{dt} = \mathbf{m}(t); \quad \frac{d\mathbf{m}(t)}{dt} = -\gamma\mathbf{m}(t) + f(\mathbf{h}(t), t, \theta),$$

satisfies the following first-order ODEs

$$\frac{d\mathbf{a}_h}{dt} = \mathbf{a}_m(t) \frac{\partial f}{\partial \mathbf{h}}; \quad \frac{d\mathbf{a}_m}{dt} = -\mathbf{a}_h(t) + \gamma\mathbf{a}_m(t). \quad (14)$$

Proof of Theorem 2.

- First, recall the adjoint state $\mathbf{a}(t)$ of the neural ODE $\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta)$ satisfy the following ODE

$$\frac{d\mathbf{a}(t)}{dt} = -\mathbf{a}(t) \frac{\partial f(\mathbf{h}(t), t, \theta)}{\partial \mathbf{h}(t)}.$$

- Let $\mathbf{n}(t) = [\mathbf{h}(t) \ \mathbf{m}(t)]$ and $\mathbf{a}(t) = [\mathbf{a}_h(t) \ \mathbf{a}_m(t)]$, then we can rewrite HBNODE as

$$\frac{d}{dt}[\mathbf{h}(t), \ \mathbf{m}(t)] = [\mathbf{m}(t), \ -\gamma\mathbf{m}(t) + f(\mathbf{h}(t), t, \theta)]$$

- Therefore, we have

$$\frac{d}{dt}[\mathbf{a}_h \ \mathbf{a}_m] = -[\mathbf{a}_h \ \mathbf{a}_m] \left[\begin{array}{cc} \frac{\partial \mathbf{m}}{\partial \mathbf{h}} & \frac{\partial \mathbf{m}}{\partial \mathbf{m}} \\ \frac{\partial(-\gamma\mathbf{m}(t)+f(\mathbf{h}(t),t,\theta))}{\partial \mathbf{h}} & \frac{\partial(-\gamma\mathbf{m}(t)+f(\mathbf{h}(t),t,\theta))}{\partial \mathbf{m}} \end{array} \right]$$

i.e.,

$$\frac{d\mathbf{a}_h}{dt} = \mathbf{a}_m(t) \frac{\partial f}{\partial \mathbf{h}}; \quad \frac{d\mathbf{a}_m}{dt} = -\mathbf{a}_h(t) + \gamma\mathbf{a}_m(t).$$

Adjoint ODE of the HBNODE

Theorem 3. The adjoint state $\mathbf{a}(t) := \partial\mathcal{L}/\partial\mathbf{h}(t)$ for the HBNODE (12) satisfies the following HBODE with the same damping parameter γ as that in (12),

$$\frac{d^2\mathbf{a}(t)}{dt^2} - \gamma\frac{d\mathbf{a}(t)}{dt} = \mathbf{a}(t)\frac{\partial f}{\partial\mathbf{h}}(\mathbf{h}(t), t, \theta). \quad (15)$$

Remark. Note that we solve the adjoint equation (15) from time $t = T$ to $t = t_0$ in the backward propagation. By letting $\tau = T - t$ and $\mathbf{b}(\tau) = \mathbf{a}(T - \tau)$, we can rewrite (15) as follows (note that $\frac{d\mathbf{a}(t)}{dt} = \frac{d\mathbf{a}(T-\tau)}{dt} = -\frac{d\mathbf{a}(T-\tau)}{d\tau} = -\frac{d\mathbf{b}(\tau)}{d\tau}$),

$$\frac{d^2\mathbf{b}(\tau)}{d\tau^2} + \gamma\frac{d\mathbf{b}(\tau)}{d\tau} = \mathbf{b}(\tau)\frac{\partial f}{\partial\mathbf{h}}(\mathbf{h}(T - \tau), T - \tau, \theta). \quad (16)$$

The adjoint of HBNODE is also a HBNODE with the same damping parameter.

Theorem 4. The adjoint states $\frac{\partial L}{\partial \mathbf{h}_t}$ and $\frac{\partial L}{\partial \mathbf{m}_t}$ of the following HBNODE

$$\frac{d\mathbf{h}(t)}{dt} = \mathbf{m}(t); \quad \frac{d\mathbf{m}(t)}{dt} = -\gamma\mathbf{m}(t) + f(\mathbf{h}(t), t, \theta). \quad (17)$$

satisfy

$$\begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \end{bmatrix} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \mathbf{h}_T} & \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \end{bmatrix} \exp \left\{ - \underbrace{\int_T^t \begin{bmatrix} 0 & \mathbf{I} \\ \frac{\partial f}{\partial \mathbf{h}} & -\gamma \mathbf{I} \end{bmatrix} ds}_{:=M} \right\}. \quad (18)$$

Proof of Theorem 4.

- Recall for neural ODE $\frac{d\mathbf{h}(t)}{dt} = f(\mathbf{h}(t), t, \theta)$, we have

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \frac{\partial \mathbf{h}_T}{\partial \mathbf{h}_t} = \frac{\partial L}{\partial \mathbf{h}_T} \exp \left\{ - \int_T^t \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}(s), s, \theta) ds \right\},$$

- Let $\mathbf{n}(t) = [\mathbf{h}(t) \ \mathbf{m}(t)]$ and $\mathbf{a}(t) = [\mathbf{a}_h(t) \ \mathbf{a}_m(t)]$, then we can rewrite HBNODE as

$$\frac{d}{dt}[\mathbf{h}(t), \ \mathbf{m}(t)] = [\mathbf{m}(t), \ -\gamma \mathbf{m}(t) + f(\mathbf{h}(t), t, \theta)]$$

- Therefore, we have

$$\begin{aligned} \left[\frac{\partial \mathcal{L}}{\partial \mathbf{h}_t} \ \frac{\partial \mathcal{L}}{\partial \mathbf{m}_t} \right] &= \left[\frac{\partial \mathcal{L}}{\partial \mathbf{h}_T} \ \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \right] \exp \left\{ - \int_T^t \left[\begin{array}{cc} \frac{\partial \mathbf{m}}{\partial \mathbf{h}} & \frac{\partial \mathbf{m}}{\partial \mathbf{m}} \\ \frac{\partial(-\gamma \mathbf{m} + f)}{\partial \mathbf{h}} & \frac{\partial(-\gamma \mathbf{m} + f)}{\partial \mathbf{m}} \end{array} \right] ds \right\} \\ &= \left[\frac{\partial \mathcal{L}}{\partial \mathbf{h}_T} \ \frac{\partial \mathcal{L}}{\partial \mathbf{m}_T} \right] \exp \left\{ - \underbrace{\int_T^t \left[\begin{array}{cc} 0 & \mathbf{I} \\ \frac{\partial f}{\partial \mathbf{h}} & -\gamma \mathbf{I} \end{array} \right] ds}_{:= \mathbf{M}} \right\}. \end{aligned}$$

Adjoint analysis

Proposition 1. The eigenvalues of the matrix $-M$ can be paired so that the sum of each pair equals $(t - T)\gamma$.

Proof of Proposition 1. Let $\mathbf{F} = \frac{1}{t-T} \int_T^t \frac{\partial f}{\partial \mathbf{h}}(\mathbf{h}(s), s, \theta) ds$ and $\mathbf{H} = \frac{1}{t-T} \mathbf{M}$, then we have the following equation

$$\mathbf{H} = \frac{1}{t-T} \mathbf{M} = \begin{bmatrix} 0 & \mathbf{I} \\ \mathbf{F} & -\gamma \mathbf{I} \end{bmatrix}. \quad (19)$$

The the characteristic polynomial of \mathbf{H} satisfies

$$\begin{aligned} ch_{\mathbf{H}}(\lambda) &= \det(\lambda \mathbf{I} - \mathbf{H}) = \det \begin{bmatrix} \lambda \mathbf{I} & -\mathbf{I} \\ -\mathbf{F} & (\lambda + \gamma) \mathbf{I} \end{bmatrix} \\ &= \det(\lambda(\lambda + \gamma) \mathbf{I} - \mathbf{F}) = -ch_{\mathbf{F}}(\lambda(\lambda + \gamma)). \end{aligned} \quad (20)$$

We can write the characteristic polynomial of \mathbf{F} as

$$ch_{\mathbf{F}}(\lambda) = \prod_{i=1}^n (\lambda - \lambda_{\mathbf{F},i}),$$

therefore

$$ch_{\mathbf{H}}(\lambda) = -ch_{\mathbf{JF}}(\lambda(\lambda + \gamma)) = -\prod_{i=1}^n (\lambda(\lambda + \gamma) - \lambda_{\mathbf{F},i}). \quad (21)$$

Therefore, the eigenvalues of \mathbf{H} appear in n pairs with each pair satisfying the quadratic equation

$$\lambda(\lambda + \gamma) - \lambda_{\mathbf{F},i} = 0. \quad (22)$$

By Vieta's formulas, the sum of these pairs are all $-\gamma$. Therefore, the eigenvalues of \mathbf{M} comes in n pairs and the sum of each pair is $-(t - T)\gamma$.

Remark on vanishing adjoint states

- By Schur decomposition, there exists an orthogonal matrix \mathbf{Q} and an upper triangular matrix \mathbf{U} , where the diagonals of \mathbf{U} are eigenvalues of \mathbf{Q} ordered by their real parts, s.t.

$$-\mathbf{M} = \mathbf{Q}\mathbf{U}\mathbf{Q}^\top \implies \exp\{-\mathbf{M}\} = \mathbf{Q} \exp\{\mathbf{U}\} \mathbf{Q}^\top. \quad (23)$$

- Let $\mathbf{v}^\top := \begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_T} & \frac{\partial L}{\partial \mathbf{m}_T} \end{bmatrix} \mathbf{Q}$,

$$\begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_t} & \frac{\partial L}{\partial \mathbf{m}_t} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_T} & \frac{\partial L}{\partial \mathbf{m}_T} \end{bmatrix} \exp\{-\mathbf{M}\} = \begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_T} & \frac{\partial L}{\partial \mathbf{m}_T} \end{bmatrix} \mathbf{Q} \exp\{\mathbf{U}\} \mathbf{Q}^\top = \mathbf{v}^\top \exp\{\mathbf{U}\} \mathbf{Q}^\top. \quad (24)$$

- By taking the ℓ_2 norm in (24) and dividing both sides by $\left\| \begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_T} & \frac{\partial L}{\partial \mathbf{m}_T} \end{bmatrix} \right\|_2$, we arrive at

$$\frac{\left\| \begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_t} & \frac{\partial L}{\partial \mathbf{m}_t} \end{bmatrix} \right\|_2}{\left\| \begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_T} & \frac{\partial L}{\partial \mathbf{m}_T} \end{bmatrix} \right\|_2} = \frac{\left\| \mathbf{v}^\top \exp\{\mathbf{U}\} \mathbf{Q}^\top \right\|_2}{\left\| \mathbf{v}^\top \mathbf{Q}^\top \right\|_2} = \frac{\left\| \mathbf{v}^\top \exp\{\mathbf{U}\} \right\|_2}{\left\| \mathbf{v} \right\|_2} = \left\| \mathbf{e}^\top \exp\{\mathbf{U}\} \right\|_2, \quad (25)$$

i.e., $\left\| \begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_t} & \frac{\partial L}{\partial \mathbf{m}_t} \end{bmatrix} \right\|_2 = \left\| \mathbf{e}^\top \exp\{\mathbf{U}\} \right\|_2 \left\| \begin{bmatrix} \frac{\partial L}{\partial \mathbf{h}_T} & \frac{\partial L}{\partial \mathbf{m}_T} \end{bmatrix} \right\|_2$ where $\mathbf{e} = \mathbf{v} / \|\mathbf{v}\|_2$.

Remark on vanishing adjoint states

- For a given constant $a > 0$, we can group the upper triangular matrix $\exp\{\mathbf{U}\}$ as follows

$$\exp\{\mathbf{U}\} := \begin{bmatrix} \exp\{\mathbf{U}_L\} & \mathbf{P} \\ 0 & \exp\{\mathbf{U}_V\} \end{bmatrix}, \quad (26)$$

where the diagonal of \mathbf{U}_L (\mathbf{U}_V) contains eigenvalues of $-\mathbf{M}$ that are no less (greater) than $(t - T)a$.

- Then, we have $\|\mathbf{e}^\top \exp\{\mathbf{U}\}\|_2 \geq \|\mathbf{e}_L^\top \exp\{\mathbf{U}_L\}\|_2$ where the vector \mathbf{e}_L denotes the first m columns of \mathbf{e} with m be the number of columns of \mathbf{U}_L . By choosing $0 \leq \gamma \leq 2a$, for every pair of eigenvalues of $-\mathbf{M}$ there is at least one eigenvalue whose real part is no less than $(t - T)a$. Therefore, $\exp\{\mathbf{U}_L\}$ decays at a rate at most $(t - T)a$, and the dimension of \mathbf{U}_L is at least $n \times n$.