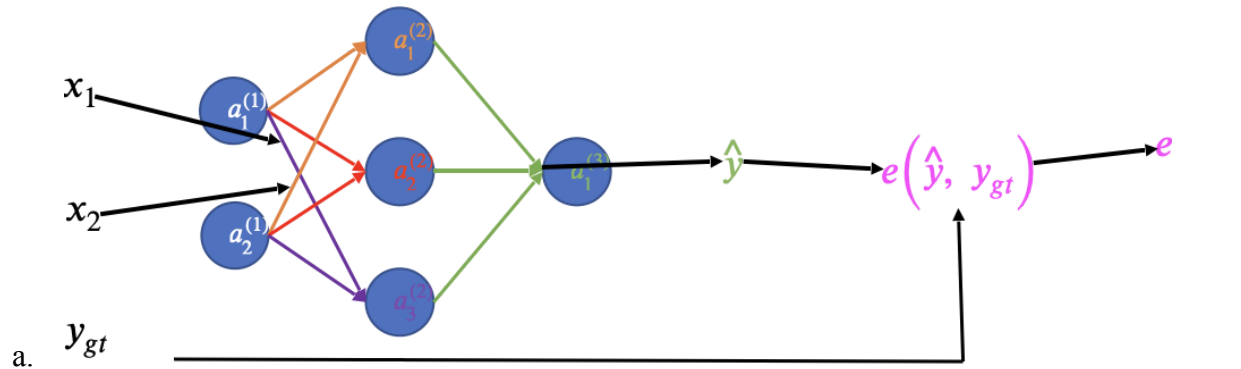


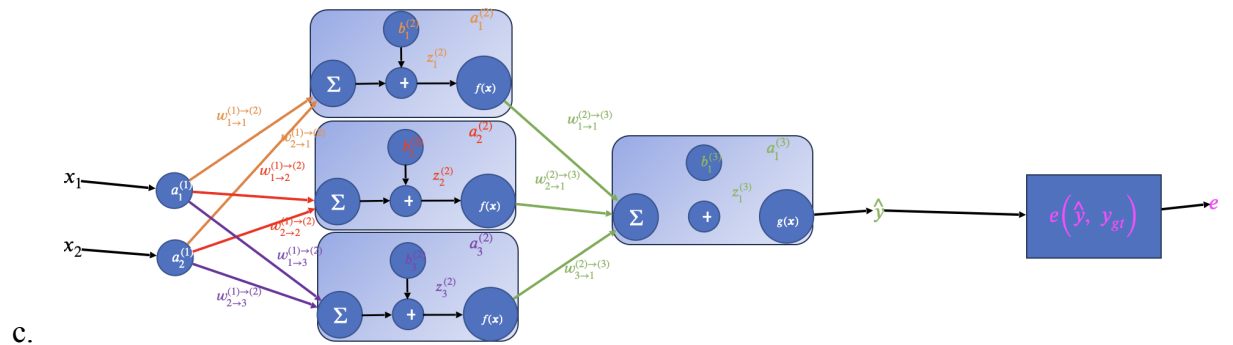
Supervised Learning VIII – Neural Networks (cont.)

1. Example



b.

$$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}} \quad g(x) = I(x) = x \quad e(\hat{y}, y_{gt}) = (y_{gt} - \hat{y})^2$$



Forward equations:

d.

$$\begin{aligned} z_1^{(2)} &= w_{1 \rightarrow 1}^{(1) \rightarrow (2)} a_1^{(1)} + w_{2 \rightarrow 1}^{(1) \rightarrow (2)} a_2^{(1)} + b_1^{(2)} & a_1^{(2)} &= f(z_1^{(2)}) \\ z_2^{(2)} &= w_{1 \rightarrow 2}^{(1) \rightarrow (2)} a_1^{(1)} + w_{2 \rightarrow 2}^{(1) \rightarrow (2)} a_2^{(1)} + b_2^{(2)} & a_2^{(2)} &= f(z_2^{(2)}) \\ z_3^{(2)} &= w_{1 \rightarrow 3}^{(1) \rightarrow (2)} a_1^{(1)} + w_{2 \rightarrow 3}^{(1) \rightarrow (2)} a_2^{(1)} + b_3^{(2)} & a_3^{(2)} &= f(z_3^{(2)}) \\ z_1^{(3)} &= w_{1 \rightarrow 1}^{(2) \rightarrow (3)} a_1^{(2)} + w_{2 \rightarrow 1}^{(2) \rightarrow (3)} a_2^{(2)} + w_{3 \rightarrow 1}^{(2) \rightarrow (3)} a_3^{(2)} + b_1^{(3)} & a_1^{(3)} &= f(z_1^{(3)}) \end{aligned}$$

Backward equations:

$$\frac{\partial e}{\partial \mathbf{w}_{2 \rightarrow 1}^{(2) \rightarrow (3)}} = \frac{\partial e}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_1^{(3)}} \frac{\partial \mathbf{z}_1^{(3)}}{\partial \mathbf{w}_{2 \rightarrow 1}^{(2) \rightarrow (3)}} = -\hat{\mathbf{y}} (y_{gt} - \hat{\mathbf{y}})^1 \mathbf{a}_2^{(2)}$$

• • •

$$\frac{\partial e}{\partial w_{1 \rightarrow 1}^{(1) \rightarrow (2)}} = \frac{\partial \textcolor{violet}{e}}{\partial \hat{\textcolor{violet}{y}}} \frac{\partial \hat{\textcolor{violet}{y}}}{\partial z_1^{(3)}} \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial w_{1 \rightarrow 1}^{(1) \rightarrow (2)}} = -\hat{\textcolor{violet}{y}} \left(y_{gt} - \hat{\textcolor{violet}{y}} \right)^1 w_{1 \rightarrow 1}^{(2) \rightarrow (3)} \sigma' \left(z_1^{(2)} \right)$$

● ● ●

$$\frac{\partial e}{\partial w_{1 \rightarrow 2}^{(1) \rightarrow (2)}} = \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(3)}}{\partial a_2^{(2)}} \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \frac{\partial z_2^{(2)}}{\partial w_{1 \rightarrow 2}^{(1) \rightarrow (2)}} = -\hat{y} \left(y_{gt} - \hat{y} \right)^1 w_{2 \rightarrow 1}^{(2) \rightarrow (3)} \sigma' \left(z_2^{(2)} \right)$$

• • •

$$\frac{\partial e}{\partial w_{2 \rightarrow 3}^{(1) \rightarrow (2)}} = \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(3)}}{\partial a_3^{(2)}} \frac{\partial a_3^{(2)}}{\partial z_3^{(2)}} \frac{\partial z_3^{(2)}}{\partial w_{2 \rightarrow 3}^{(1) \rightarrow (2)}} = -\hat{y} \left(y_{gt} - \hat{y} \right)^1 w_{2 \rightarrow 1}^{(2) \rightarrow (3)} \sigma' \left(z_3^{(2)} \right)$$

• • •

$$\frac{\partial \mathbf{e}}{\partial \mathbf{b}_1^{(3)}} = \frac{\partial \mathbf{e}}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_1^{(3)}} \frac{\partial \mathbf{z}_1^{(3)}}{\partial \mathbf{b}_1^{(3)}} = -\hat{\mathbf{y}}(y_{gt} - \hat{\mathbf{y}})^1 \mathbf{1}$$

$$\frac{\partial e}{\partial b_1^{(2)}} = \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial b_1^{(2)}} = -\hat{y} (y_{gt} - \hat{y})^1 w_{1 \rightarrow 1}^{(2) \rightarrow (3)} \sigma'(z_1^{(2)})$$

Backward equations:

$$\frac{\partial e}{\partial w_{2 \rightarrow 1}^{(2) \rightarrow (3)}} = \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(3)}}{\partial w_{2 \rightarrow 1}^{(2) \rightarrow (3)}} \dots$$

/ other

$$\frac{\partial e}{\partial w_{i \rightarrow 1}^{(1) \rightarrow (2)}} = \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \underbrace{\left(\frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right)}_{\text{Shared w/ other}} \frac{\partial z_1^{(2)}}{\partial w_{i \rightarrow 1}^{(1) \rightarrow (2)}} \frac{\partial e}{\partial w_{i \rightarrow j}^{(1) \rightarrow (2)}} \bullet \bullet \bullet$$

~~Shared w/ other~~

$$\frac{\partial e}{\partial w_{1 \rightarrow 2}^{(1) \rightarrow (2)}} = \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(3)}}{\partial a_2^{(2)}} \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \frac{\partial z_2^{(2)}}{\partial w_{1 \rightarrow 2}^{(1) \rightarrow (2)}} \frac{\partial e}{\partial w_{i \rightarrow j}^{(1) \rightarrow (2)}} \quad \text{Shared w/ other}$$

$$\frac{\partial \mathbf{e}}{\partial \mathbf{w}_{2 \rightarrow 3}^{(1 \rightarrow 2)}} = \frac{\partial \mathbf{e}}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{z}_1^{(3)}} \left(\frac{\sigma \mathbf{z}_1^{(3)}}{\partial \mathbf{a}_3^{(2)}} \frac{\partial \mathbf{a}_3^{(2)}}{\partial \mathbf{z}_3^{(2)}} \frac{\partial \mathbf{z}_3^{(2)}}{\partial \mathbf{w}_{2 \rightarrow 3}^{(1 \rightarrow 2)}} \right) \bullet \bullet \bullet$$

Only compute once!

Shared w/ other

$$\frac{\partial \mathbf{e}}{\partial \mathbf{w}_{i \rightarrow j}^{(1) \rightarrow (2)}}$$

Only compute once!

$$\frac{\partial e}{\partial b_1^{(3)}} = \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(3)}}{\partial b_1^{(3)}} + \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(2)}}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \frac{\partial z_1^{(2)}}{\partial b_1^{(2)}} + \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(2)}}{\partial a_2^{(2)}} \frac{\partial a_2^{(2)}}{\partial z_2^{(2)}} \frac{\partial z_2^{(2)}}{\partial b_2^{(2)}} + \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z_1^{(3)}} \frac{\partial z_1^{(2)}}{\partial a_3^{(2)}} \frac{\partial a_3^{(2)}}{\partial z_3^{(2)}} \frac{\partial z_3^{(2)}}{\partial b_3^{(2)}}$$

f.

- Cache the value and compute it only once and store them

$$\frac{\partial e}{\partial \mathbf{z}_1^{(3)}} = \frac{\partial e}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial \mathbf{z}_1^{(3)}} = -\hat{y}(y_{gt} - \hat{y})g'(\mathbf{z}_1^{(3)})$$

$$\frac{\partial e}{\partial \mathbf{z}_1^{(2)}} = \frac{\partial e}{\partial \mathbf{z}_1^{(3)}} \frac{\partial \mathbf{z}_1^{(3)}}{\partial \mathbf{a}_1^{(2)}} \frac{\partial \mathbf{a}_1^{(2)}}{\partial \mathbf{z}_1^{(2)}} = \frac{\partial e}{\partial \mathbf{z}_1^{(3)}} w_{1 \rightarrow 1}^{(2) \rightarrow (3)} \sigma'(\mathbf{z}_1^{(2)})$$

$$\frac{\partial e}{\partial \mathbf{z}_2^{(2)}} = \frac{\partial e}{\partial \mathbf{z}_1^{(3)}} \frac{\partial \mathbf{z}_1^{(3)}}{\partial \mathbf{a}_2^{(2)}} \frac{\partial \mathbf{a}_2^{(2)}}{\partial \mathbf{z}_2^{(2)}} = \frac{\partial e}{\partial \mathbf{z}_1^{(3)}} w_{2 \rightarrow 1}^{(2) \rightarrow (3)} \sigma'(\mathbf{z}_2^{(2)})$$

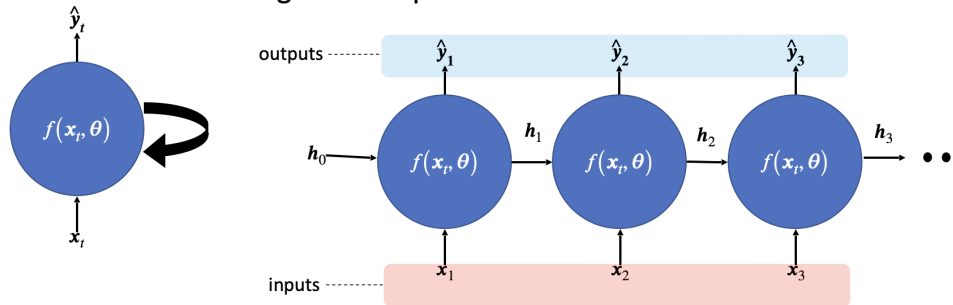
$$\text{g. } \frac{\partial e}{\partial \mathbf{z}_3^{(2)}} = \frac{\partial e}{\partial \mathbf{z}_1^{(3)}} \frac{\partial \mathbf{z}_1^{(3)}}{\partial \mathbf{a}_3^{(2)}} \frac{\partial \mathbf{a}_3^{(2)}}{\partial \mathbf{z}_3^{(2)}} = \frac{\partial e}{\partial \mathbf{z}_1^{(3)}} w_{3 \rightarrow 1}^{(2) \rightarrow (3)} \sigma'(\mathbf{z}_3^{(2)})$$

i. Cached values

2. Backpropagation Through Time (BPTT)

- What about sequences? Current NNs only deal with grids/vector input
- Recurrent NNs (RNNs): derived from State Machines
- “Unrolls” given a sequence

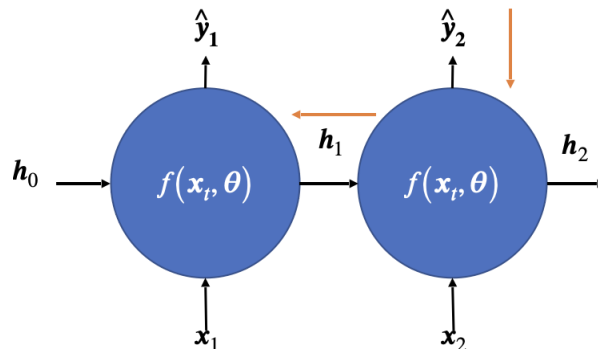
“Unrolls” given a sequence:



i.

3. Vanilla RNN with BPTT

- Using computation and perform backprop like normal



$$\mathbf{H}_t = f(\mathbf{X}_t \mathbf{U} + \mathbf{H}_{t-1} \mathbf{W} + \mathbf{b}_s)$$

$$\mathbf{Y}_t = \mathbf{H}_t \mathbf{V} + \mathbf{b}_o$$

i.

1. Input is in batch form

$$\nabla_{\theta} L = \begin{matrix} \frac{\partial L}{\partial \mathbf{b}_o} \\ \frac{\partial L}{\partial \mathbf{b}_s} \\ \frac{\partial L}{\partial \mathbf{U}} \\ \frac{\partial L}{\partial \mathbf{W}} \\ \frac{\partial L}{\partial \mathbf{V}} \end{matrix} = \begin{matrix} \frac{\partial L}{\partial \mathbf{b}_o} | x_2, y_2 \\ \frac{\partial L}{\partial \mathbf{b}_s} | x_2, y_2 \\ \frac{\partial L}{\partial \mathbf{U}} | x_2, y_2 \\ \frac{\partial L}{\partial \mathbf{W}} | x_2, y_2 \\ \frac{\partial L}{\partial \mathbf{V}} | x_2, y_2 \end{matrix} + \begin{matrix} \frac{\partial L}{\partial \mathbf{b}_o} | x_1, y_1 \\ \frac{\partial L}{\partial \mathbf{b}_s} | x_1, y_1 \\ \frac{\partial L}{\partial \mathbf{U}} | x_1, y_1 \\ \frac{\partial L}{\partial \mathbf{W}} | x_1, y_1 \\ \frac{\partial L}{\partial \mathbf{V}} | x_1, y_1 \end{matrix}$$

ii.

4. RNN Variants

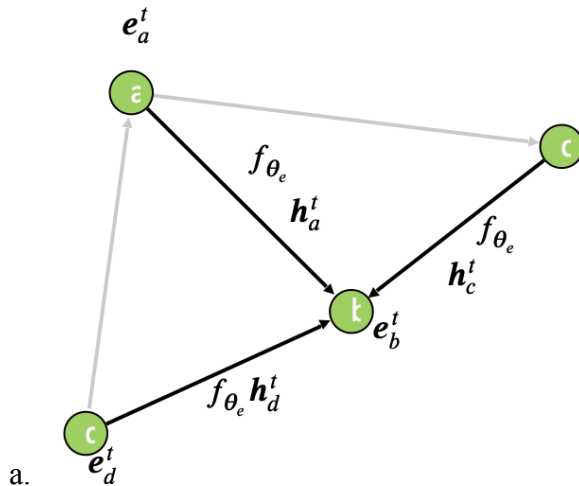
- a. Different “cell types”

- i. LSTM
- ii. GRU
- iii. LRU
- iv. ...

- b. Only difference is in how internal state is computed/maintained

- i. Big difference, add specific areas for “forgetting” and “remembering”

5. Graph Structure?



What you need to define:

An **edge function** f_{θ_e}

An **aggregator function** γ

A **vertex function** g_{θ_v}

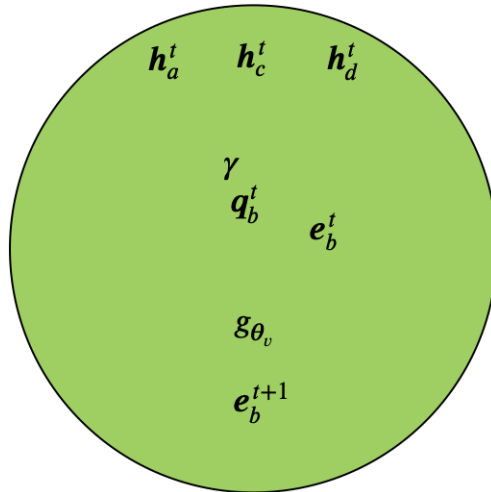
Computational Process:

$$1) h_i^t = f_{\theta_e}(e_i^t), i \in N(b)$$

$$2) q_b^t = \gamma(\{h_i^t\}_{i \in N(b)})$$

$$3) e_b^{t+1} = g_{\theta_v}(e_b^t, q_b^t)$$

b.



- c.
6. Easy to Vectorize

a. $\gamma = \sum$ (matrix mult. w/ adjacency matrix)

b. $f_{\theta_e} = \text{dense NN}$

c. $g_{\theta_v} =$ some (element-wise independent) activation function 'g'

next layer

f_{θ_e}

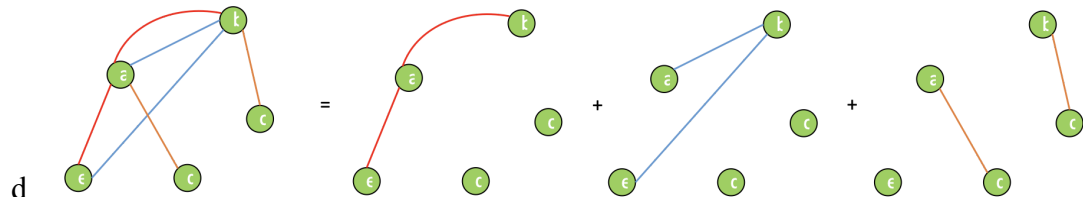
$$E^{t+1} = g(A^t E^t W^t)$$

vectors for all vertices at current layer

adjacency matrix for current layer

- d.
7. Multi-view Graphs?

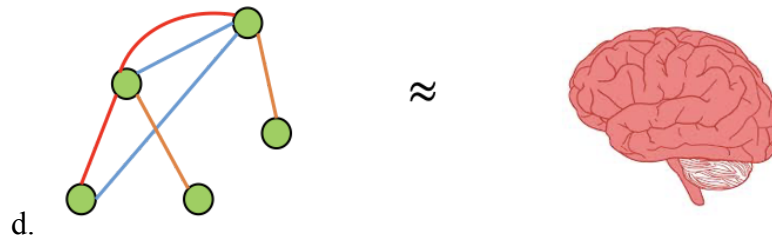
- a. Edges all have one "type" (represent the same type of link)
- b. Not true for many problems
- c. Can we make GNNs work for multi-view graphs?



$$E^{t+1} = g\left(\frac{1}{3} A^t E^t W^t + \frac{1}{3} B^t E^t W^t + \frac{1}{3} C^t E^t W^t\right)$$

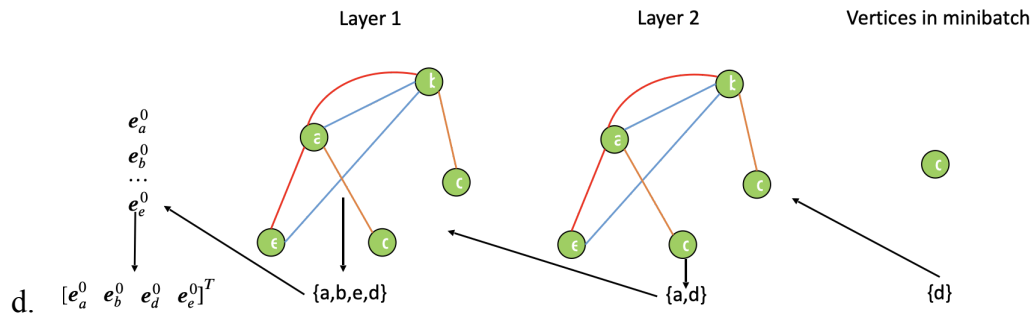
- e.
8. Can we use GNNs to model the brain?

- a. Built in topology
- b. Encode lifecycle rules
 - i. Start off with lots of connections
 - ii. Prune connections over time
- c. Natural architecture for wiring costs, energy constraints, etc.



9. Make GNNs Computable

- a. $\mathbf{E}^{t+1} = g(\mathbf{A}^t \mathbf{E}^t \mathbf{W}^t)$ takes a lot of memory (and time)
- b. Problem: too big for (most) GPUs....doesn't scale
 - i. Limitations in sparse support
- c. Can we make it?



10. Changing the Graph Structure: Vertex Inflation

- a. The brain changes its structure over its lifetime
- b. GNNs reuse the same graph for each layer
- c. Idea: change the graph structure as a function of layer
 - i. Vertex inflation: shortcut “long paths” in the graph
 - ii. Performance?

