

Tutorial: Neural Networks & Weighted Automata

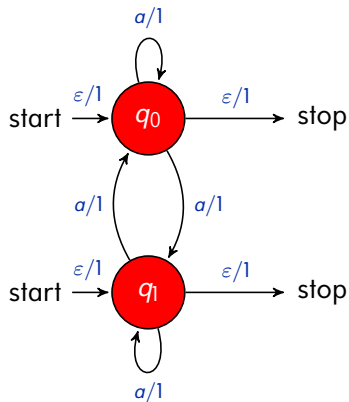
Andreas Maletti

Universität Leipzig, Germany

WATA — April 20, 2021

Weighted automata

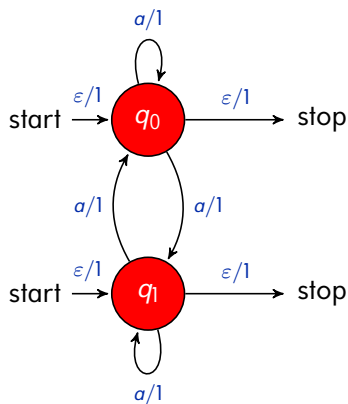
Weight structure: commutative field $(\mathbb{R}, +, \cdot, 0, 1)$



Path from start to stop: Product over labels in $\Sigma^* \times \mathbb{R}$

Weighted automata

- 2^{k+1} paths for input a^k
- Every path has weight 1

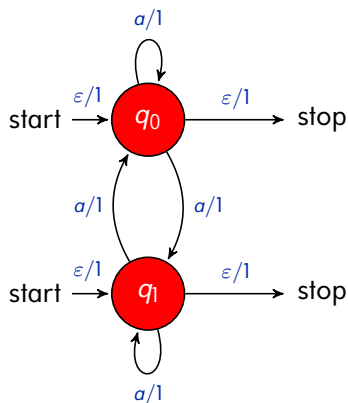


Weighted automata

- 2^{k+1} paths for input a^k
- Every path has weight 1

Weight of input:

$$\text{wt}(a^k) = \sum_{\pi \text{ path for } a^k} \text{wt}(\pi) = 2^{k+1}$$



Definition (weighted automaton)

Weighted automaton $\mathcal{A} = (Q, \Sigma, I, (M_a)_{a \in \Sigma}, F)$

- finite set Q
- alphabet Σ
- $I \in \mathbb{R}^Q$ initial weight vector
- $M_a \in \mathbb{R}^{Q \times Q}$ transition weight matrix for each $a \in \Sigma$
- $F \in \mathbb{R}^Q$ final weight vector

Weighted automata

Definition (weighted automaton)

Weighted automaton $\mathcal{A} = (Q, \Sigma, I, (M_a)_{a \in \Sigma}, F)$

- finite set Q
- alphabet Σ
- $I \in \mathbb{R}^Q$ initial weight vector
- $M_a \in \mathbb{R}^{Q \times Q}$ transition weight matrix for each $a \in \Sigma$
- $F \in \mathbb{R}^Q$ final weight vector

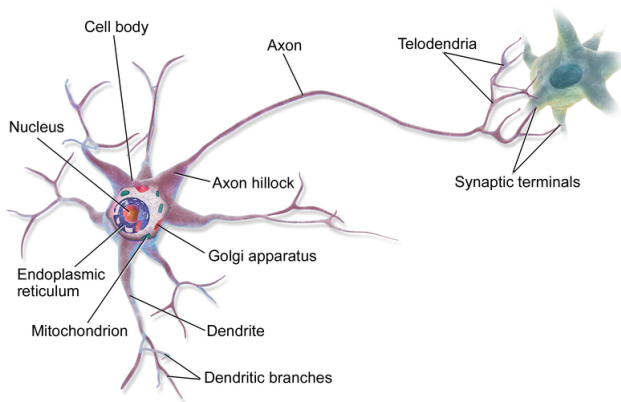
Semantics

$$\mathcal{A}(a_1 \cdots a_k) = I^T M_{a_1} \cdots M_{a_k} F \quad \text{for all } a_1, \dots, a_k \in \Sigma$$

Neural networks

Neural network

- Edge-weighted graph (N, E, wt) of neurons inspired by nerve cells
- Focus on Recurrent Neural Nets (RNN)



Neural networks

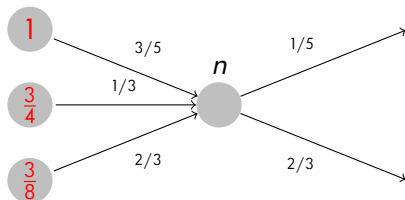
Formalization

- Each neuron $n \in N$ is vertex
- Weight $wt(n', n) \in \mathbb{R}$ indicates connection strength from n' to n
- Incoming edges of n contribute to activation of n
(All input potentials moderated by connection strength added up)

Neural networks

Formalization

- Each neuron $n \in N$ is vertex
- Weight $wt(n', n) \in \mathbb{R}$ indicates connection strength from n' to n
- Incoming edges of n contribute to activation of n
(All input potentials moderated by connection strength added up)



Input potential

$$1 \cdot \frac{3}{5} + \frac{3}{4} \cdot \frac{1}{3} + \frac{3}{8} \cdot \frac{2}{3} = \frac{3}{5} + \frac{1}{2} = 1.1$$

Neural networks

Formalization

- Given state $s: N \rightarrow \mathbb{R}$ and **activation function** $\sigma: \mathbb{R} \rightarrow \mathbb{R}$
- Input potential supplied to σ yields new state $s': N \rightarrow \mathbb{R}$

$$s' = \sigma(\text{wt} \cdot s)$$

$$s'(n) = \sigma \left\langle \sum_{(n',n) \in E} \text{wt}(n',n) \cdot s(n') \right\rangle$$

for all $n \in N$

Neural networks

Formalization

- Given state $s: N \rightarrow \mathbb{R}$ and **activation function** $\sigma: \mathbb{R} \rightarrow \mathbb{R}$
- Input potential supplied to σ yields new state $s': N \rightarrow \mathbb{R}$

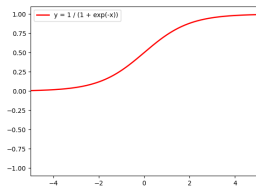
$$s' = \sigma(\text{wt} \cdot s)$$

$$s'(n) = \sigma \left\langle \sum_{(n',n) \in E} \text{wt}(n',n) \cdot s(n') \right\rangle$$

for all $n \in N$

Standard activation functions

Sigmoid



Neural networks

Formalization

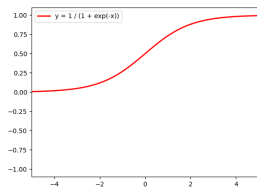
- Given state $s: N \rightarrow \mathbb{R}$ and **activation function** $\sigma: \mathbb{R} \rightarrow \mathbb{R}$
- Input potential supplied to σ yields new state $s': N \rightarrow \mathbb{R}$

$$s' = \sigma \langle \text{wt} \cdot s \rangle \qquad s'(n) = \sigma \left\langle \sum_{(n',n) \in E} \text{wt}(n',n) \cdot s(n') \right\rangle$$

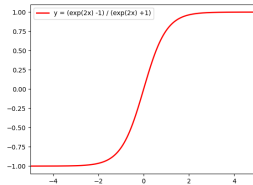
for all $n \in N$

Standard activation functions

Sigmoid



Tanh



Neural networks

Formalization

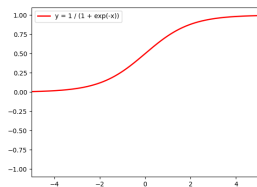
- Given state $s: N \rightarrow \mathbb{R}$ and **activation function** $\sigma: \mathbb{R} \rightarrow \mathbb{R}$
- Input potential supplied to σ yields new state $s': N \rightarrow \mathbb{R}$

$$s' = \sigma \langle \text{wt} \cdot s \rangle \qquad s'(n) = \sigma \left\langle \sum_{(n',n) \in E} \text{wt}(n',n) \cdot s(n') \right\rangle$$

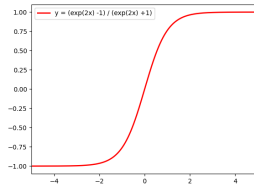
for all $n \in N$

Standard activation functions

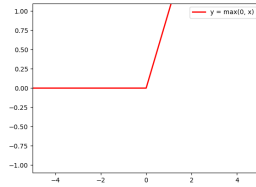
Sigmoid



Tanh

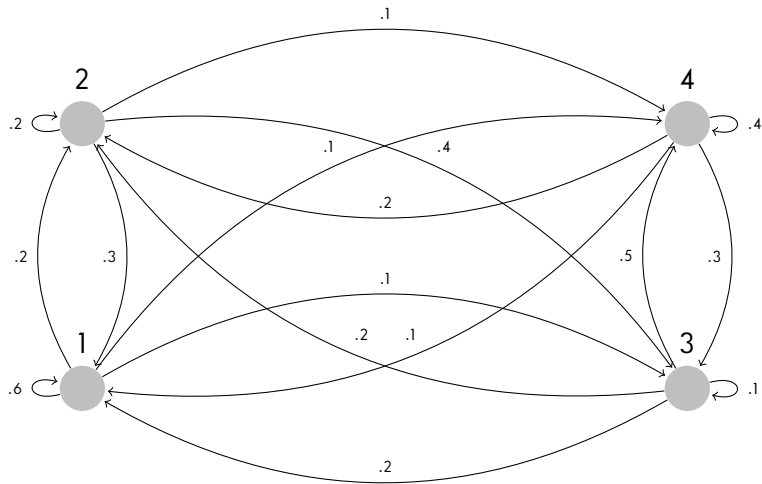


ReLU



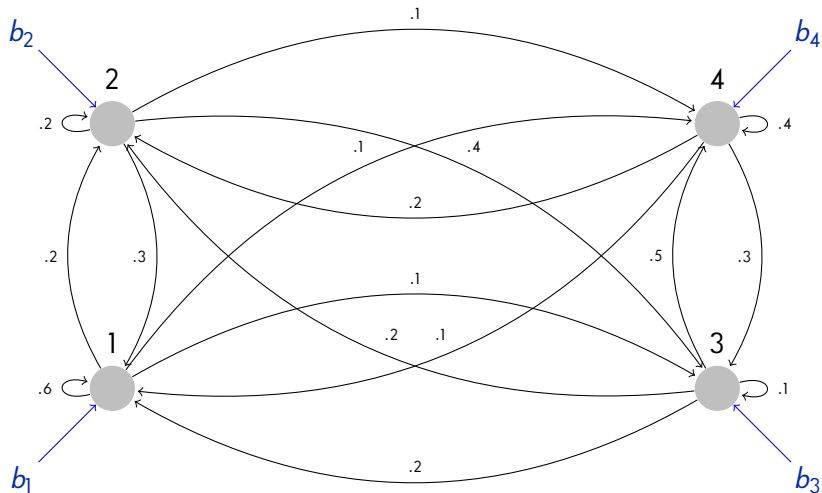
Neural networks

RNN = edge-labeled complete graph $(N, N \times N, wt)$ with $wt \in \mathbb{R}^{N \times N}$



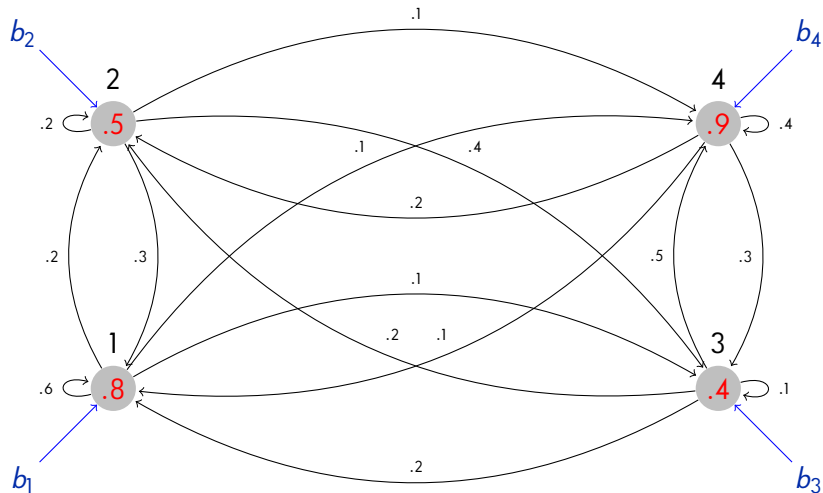
Neural networks

RNN = edge-labeled complete graph $(N, N \times N, wt)$ with $wt \in \mathbb{R}^{N \times N}$



Neural networks

RNN = edge-labeled complete graph $(N, N \times N, wt)$ with $wt \in \mathbb{R}^{N \times N}$



State behavior

- Edge weight matrix $\text{wt} \in \mathbb{R}^{N \times N}$ and input vector $\text{inp} \in \mathbb{R}^N$

$$\text{wt} = \begin{pmatrix} .6 & .3 & .2 & .1 \\ .2 & .2 & .2 & .2 \\ .1 & .4 & .1 & .3 \\ .1 & .1 & .5 & .4 \end{pmatrix} \quad \text{inp} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{pmatrix}$$

State behavior

- Edge weight matrix $\text{wt} \in \mathbb{R}^{N \times N}$ and input vector $\text{inp} \in \mathbb{R}^N$

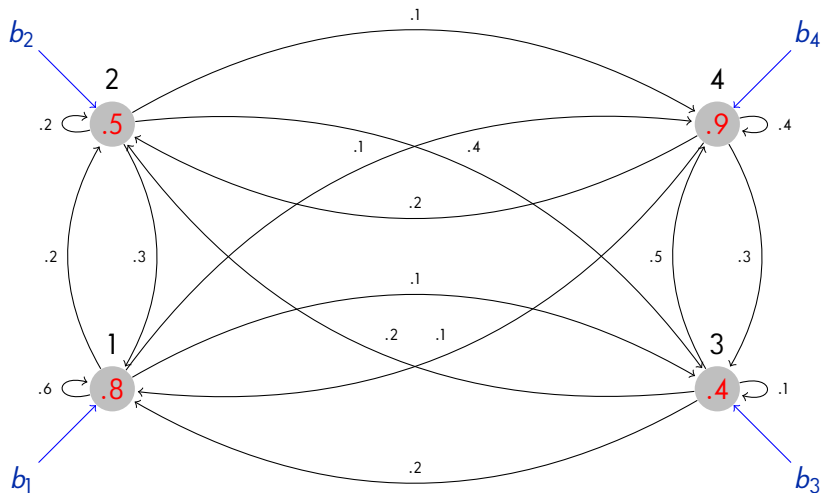
$$\text{wt} = \begin{pmatrix} .6 & .3 & .2 & .1 \\ .2 & .2 & .2 & .2 \\ .1 & .4 & .1 & .3 \\ .1 & .1 & .5 & .4 \end{pmatrix} \qquad \text{inp} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{pmatrix}$$

- State vector $s \in \mathbb{R}^N$ and ReLU activation $\sigma(x) = \max(0, x)$
- Follow-state vector $s' \in \mathbb{R}^N$

$$s' = \sigma(\text{inp} + \text{wt} \cdot s) \qquad s'(n) = \sigma\left(\text{inp}(n) + \sum_{n' \in N} \text{wt}(n', n) \cdot s(n')\right)$$

for all $n \in N$

Neural networks



Example

$$\text{wt} = \begin{pmatrix} .6 & .3 & .2 & .1 \\ .2 & .2 & .2 & .2 \\ .1 & .4 & .1 & .3 \\ .1 & .1 & .5 & .4 \end{pmatrix}$$

$$s = \begin{pmatrix} .25 \\ .25 \\ .25 \\ .25 \end{pmatrix}$$

$$\text{inp} = \begin{pmatrix} -.8 \\ .5 \\ .4 \\ -.3 \end{pmatrix}$$

Neural networks

Example

$$\text{wt} = \begin{pmatrix} .6 & .3 & .2 & .1 \\ .2 & .2 & .2 & .2 \\ .1 & .4 & .1 & .3 \\ .1 & .1 & .5 & .4 \end{pmatrix} \quad s = \begin{pmatrix} .25 \\ .25 \\ .25 \\ .25 \end{pmatrix} \quad \text{inp} = \begin{pmatrix} -.8 \\ .5 \\ .4 \\ -.3 \end{pmatrix}$$

Then

$$\text{wt} \cdot s = \begin{pmatrix} .3 \\ .2 \\ .225 \\ .275 \end{pmatrix}$$

Neural networks

Example

$$\text{wt} = \begin{pmatrix} .6 & .3 & .2 & .1 \\ .2 & .2 & .2 & .2 \\ .1 & .4 & .1 & .3 \\ .1 & .1 & .5 & .4 \end{pmatrix}$$

$$s = \begin{pmatrix} .25 \\ .25 \\ .25 \\ .25 \end{pmatrix}$$

$$\text{inp} = \begin{pmatrix} -.8 \\ .5 \\ .4 \\ -.3 \end{pmatrix}$$

Then

$$\text{wt} \cdot s = \begin{pmatrix} .3 \\ .2 \\ .225 \\ .275 \end{pmatrix}$$

$$\text{inp} + \text{wt} \cdot s = \begin{pmatrix} -.5 \\ .7 \\ .625 \\ -.025 \end{pmatrix}$$

Neural networks

Example

$$\text{wt} = \begin{pmatrix} .6 & .3 & .2 & .1 \\ .2 & .2 & .2 & .2 \\ .1 & .4 & .1 & .3 \\ .1 & .1 & .5 & .4 \end{pmatrix} \quad s = \begin{pmatrix} .25 \\ .25 \\ .25 \\ .25 \end{pmatrix} \quad \text{inp} = \begin{pmatrix} -.8 \\ .5 \\ .4 \\ -.3 \end{pmatrix}$$

Then

$$\text{wt} \cdot s = \begin{pmatrix} .3 \\ .2 \\ .225 \\ .275 \end{pmatrix} \quad \text{inp} + \text{wt} \cdot s = \begin{pmatrix} -.5 \\ .7 \\ .625 \\ -.025 \end{pmatrix}$$

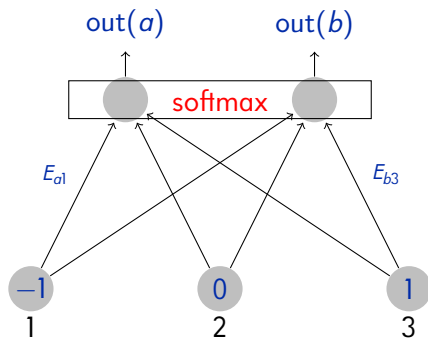
Applying ReLU activation σ

$$s' = \sigma(\text{inp} + \text{wt} \cdot s) = (0 \quad .7 \quad .625 \quad 0)^T$$

Neural networks

Adding output

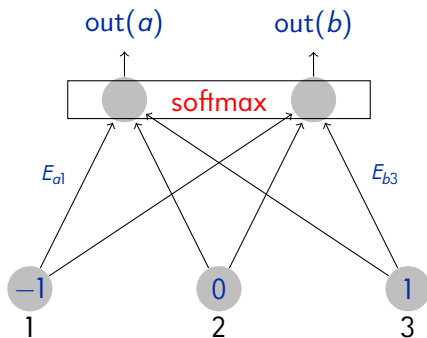
- Prediction matrix $E \in \mathbb{R}^{\Sigma \times N}$ to produce vector entry for each letter



Adding output

- Prediction matrix $E \in \mathbb{R}^{\Sigma \times N}$ to produce vector entry for each letter
- Turn into probability via softmax: $\mathbb{R}^{\Sigma} \rightarrow \mathbb{R}^{\Sigma}$ for all $v \in \mathbb{R}^{\Sigma}$ and $a \in \Sigma$

$$\text{softmax}(v)_a = \frac{e^{v_a}}{\sum_{b \in \Sigma} e^{v_b}}$$

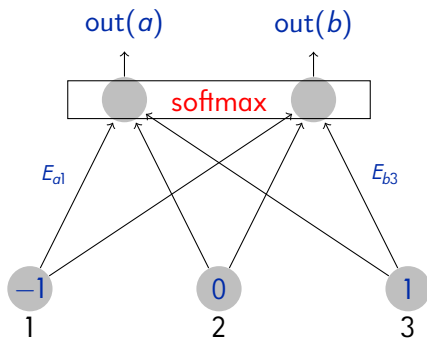


Adding output

- Prediction matrix $E \in \mathbb{R}^{\Sigma \times N}$ to produce vector entry for each letter
- Turn into probability via softmax: $\mathbb{R}^{\Sigma} \rightarrow \mathbb{R}^{\Sigma}$ for all $v \in \mathbb{R}^{\Sigma}$ and $a \in \Sigma$

$$\text{softmax}(v)_a = \frac{e^{v_a}}{\sum_{b \in \Sigma} e^{v_b}}$$

- Letter classification $\text{out} = \text{softmax}\langle E \cdot s \rangle$

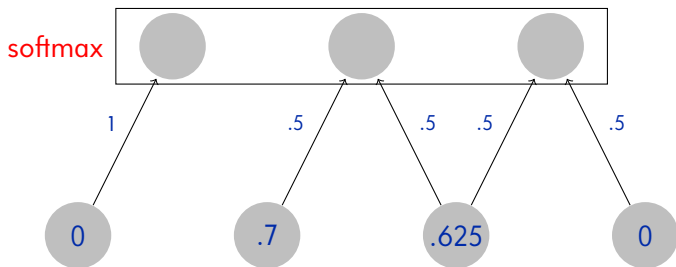


Neural networks

Example

- State $s = (0 \ .7 \ .625 \ 0)^T$ and prediction $E = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & .5 & .5 & 0 \\ 0 & 0 & .5 & .5 \end{pmatrix}$

$$\text{softmax}\langle E \cdot s \rangle = \text{softmax}\left\langle \begin{matrix} 0 \\ .6625 \\ .3125 \end{matrix} \right\rangle \approx \begin{pmatrix} .23 \\ .45 \\ .32 \end{pmatrix}$$



Definition (RNN)

(Single-layer) **RNN** is $(N, \Sigma, s_{-1}, wt, (inp_\sigma)_{\sigma \in \Sigma_\$}, E)$

- finite set N and alphabet Σ
- initial state $s_{-1} \in \mathbb{Q}^N$ and transition weights $wt \in \mathbb{Q}^{N \times N}$
- input vector $inp_\sigma \in \mathbb{Q}^N$ for every $\sigma \in \Sigma_\$$ with $\Sigma_\$ = \Sigma \cup \{\$\}$
- prediction matrix $E \in \mathbb{Q}^{\Sigma_\$ \times N}$

Definition (Semantics of RNN)

Let $\mathcal{R} = (N, \Sigma, s_{-1}, \text{wt}, (\text{inp}_\sigma)_{\sigma \in \Sigma_\$}, E)$ be RNN and $w = a_1 \cdots a_k \in \Sigma^*$

- $a_0 = \$$, $a_{k+1} = \$$, and $s_{-1}^w = s_{-1}$
- State $s_t^w = \sigma(\text{inp}_{a_t} + \text{wt} \cdot s_{t-1}^w)$ for every $0 \leq t \leq k$

Definition (Semantics of RNN)

Let $\mathcal{R} = (N, \Sigma, s_{-1}, \text{wt}, (\text{inp}_\sigma)_{\sigma \in \Sigma_s}, E)$ be RNN and $w = a_1 \cdots a_k \in \Sigma^*$

- $a_0 = \$$, $a_{k+1} = \$$, and $s_{-1}^w = s_{-1}$
- State $s_t^w = \sigma(\text{inp}_{a_t} + \text{wt} \cdot s_{t-1}^w)$ for every $0 \leq t \leq k$
- Output $\text{out}_t^w = \text{softmax}(E \cdot s_t^w)$ for every $0 \leq t \leq k$

Definition (Semantics of RNN)

Let $\mathcal{R} = (N, \Sigma, s_{-1}, wt, (\text{inp}_\sigma)_{\sigma \in \Sigma_\mathcal{S}}, E)$ be RNN and $w = a_1 \cdots a_k \in \Sigma^*$

- $a_0 = \$$, $a_{k+1} = \$$, and $s_{-1}^w = s_{-1}$
- State $s_t^w = \sigma(\text{inp}_{a_t} + wt \cdot s_{t-1}^w)$ for every $0 \leq t \leq k$
- Output $\text{out}_t^w = \text{softmax}(E \cdot s_t^w)$ for every $0 \leq t \leq k$
- Weight of w

$$\mathcal{R}(w) = \prod_{t=0}^k \text{out}_t^w(a_{t+1})$$

Definition (Semantics of RNN)

Let $\mathcal{R} = (N, \Sigma, s_{-1}, wt, (\text{inp}_\sigma)_{\sigma \in \Sigma_s}, E)$ be RNN and $w = a_1 \cdots a_k \in \Sigma^*$

- $a_0 = \$$, $a_{k+1} = \$$, and $s_{-1}^w = s_{-1}$
- State $s_t^w = \sigma \langle \text{inp}_{a_t} + wt \cdot s_{t-1}^w \rangle$ for every $0 \leq t \leq k$
- Output $\text{out}_t^w = \text{softmax} \langle E \cdot s_t^w \rangle$ for every $0 \leq t \leq k$
- Weight of w

$$\mathcal{R}(w) = \prod_{t=0}^k \text{out}_t^w(a_{t+1})$$

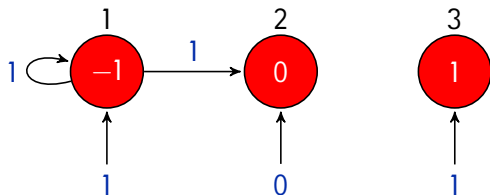
Note

- Thus $\mathcal{R}: \Sigma^* \rightarrow [0, 1]$ with $\sum_{w \in \Sigma^*} \mathcal{R}(w) \leq 1$

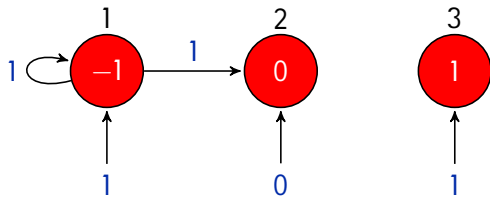
Neural networks

Example RNN $\mathcal{R} = (\{1, 2, 3\}, \{a\}, s_{-1}, wt, (inp_s, inp_a), E)$ and $z \in \mathbb{Q}$

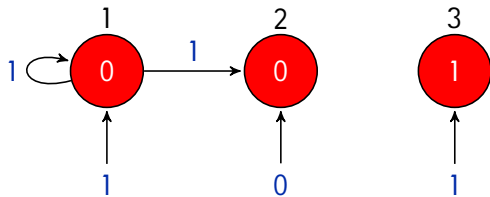
$$wt = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \quad s_{-1} = \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix} \quad inp_s = inp_a = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$



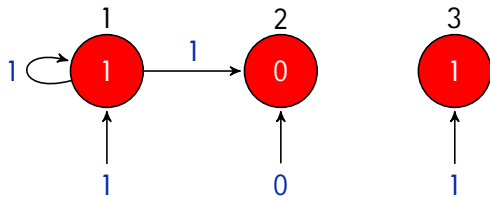
Example (cont'd)



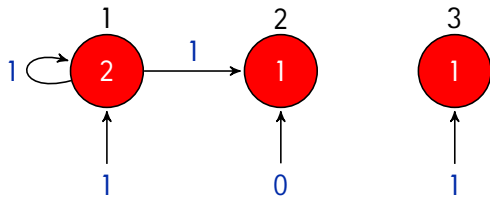
Example (cont'd)



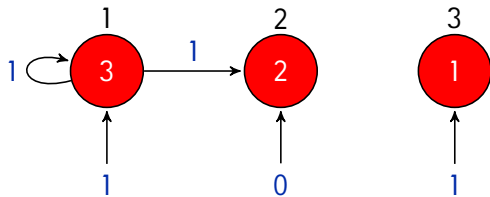
Example (cont'd)



Example (cont'd)

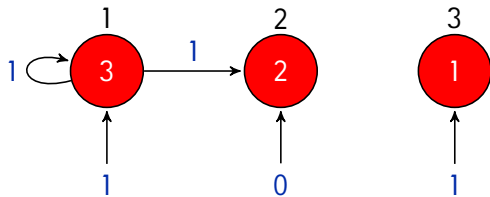


Example (cont'd)



- Input $w = a^k$

Example (cont'd)



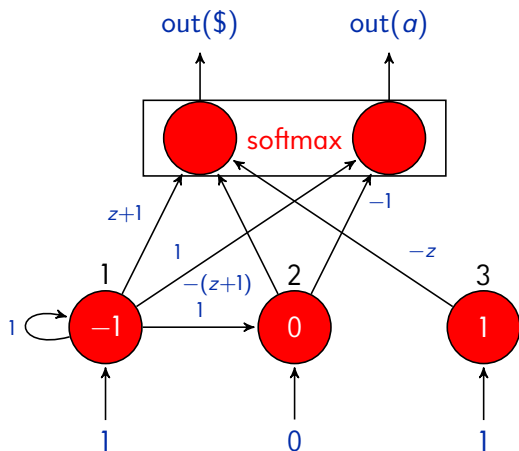
- Input $w = a^k$

- States $s_t^w = \sigma \left\langle \begin{matrix} 1 + s_{t-1}^w(1) \\ s_{t-1}^w(1) \\ 1 \end{matrix} \right\rangle = \sigma \left\langle \begin{matrix} t \\ t-1 \\ 1 \end{matrix} \right\rangle$ for all $0 \leq t \leq k$

Neural networks

Example (cont'd)

- Prediction matrix $E = \begin{pmatrix} z+1 & -(z+1) & -z \\ 1 & -1 & 0 \end{pmatrix}$



Example (cont'd)

- States $s_j^w = \sigma \left\langle \begin{matrix} t \\ t-1 \\ 1 \end{matrix} \right\rangle$ for all $0 \leq t \leq k$

Example (cont'd)

- States $s_t^w = \sigma \left\langle \begin{matrix} t \\ t-1 \\ 1 \end{matrix} \right\rangle$ for all $0 \leq t \leq k$

- Prediction matrix $E = \begin{pmatrix} z+1 & -(z+1) & -z \\ 1 & -1 & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

Example (cont'd)

- States $s_t^w = \sigma \left\langle \begin{matrix} t \\ t-1 \\ 1 \end{matrix} \right\rangle$ for all $0 \leq t \leq k$
- Prediction matrix $E = \begin{pmatrix} z+1 & -(z+1) & -z \\ 1 & -1 & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$
- $\text{out}_t^w = \text{softmax} \left\langle \begin{matrix} (z+1)t - (z+1)\sigma\langle t-1 \rangle - z \\ t - \sigma\langle t-1 \rangle \end{matrix} \right\rangle$ for all $0 \leq t \leq k$

Example (cont'd)

- States $s_t^w = \sigma \left\langle \begin{matrix} t \\ t-1 \\ 1 \end{matrix} \right\rangle$ for all $0 \leq t \leq k$
- Prediction matrix $E = \begin{pmatrix} z+1 & -(z+1) & -z \\ 1 & -1 & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$
- $\text{out}_t^w = \text{softmax} \left\langle \begin{matrix} (z+1)t - (z+1)\sigma \langle t-1 \rangle - z \\ t - \sigma \langle t-1 \rangle \end{matrix} \right\rangle$ for all $0 \leq t \leq k$
- $\text{out}_0^w = \text{softmax} \left\langle \begin{matrix} -z \\ 0 \end{matrix} \right\rangle$ and $\text{out}_t^w = \text{softmax} \left\langle \begin{matrix} 1 \\ 1 \end{matrix} \right\rangle$ for all $1 \leq t \leq k$

Example (cont'd)

- $\text{out}_0^w = \text{softmax}\left\langle \begin{matrix} -z \\ 0 \end{matrix} \right\rangle$ and $\text{out}_t^w = \text{softmax}\left\langle \begin{matrix} 1 \\ 1 \end{matrix} \right\rangle$ for all $1 \leq t \leq k$

Example (cont'd)

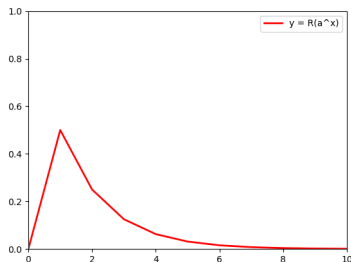
- $\text{out}_0^w = \text{softmax}\left\langle \begin{matrix} -z \\ 0 \end{matrix} \right\rangle$ and $\text{out}_t^w = \text{softmax}\left\langle \begin{matrix} 1 \\ 1 \end{matrix} \right\rangle$ for all $1 \leq t \leq k$

- Computed weights

$$\mathcal{R}(\varepsilon) = \frac{e^{-z}}{1 + e^{-z}} = \frac{1}{1 + e^z}$$

$$\mathcal{R}(\alpha^k) = \frac{1}{1 + e^{-z}} \cdot 2^{-k}$$

for $k \neq 0$



(graph for $z = 100$)

Neural networks

Example RNN $\mathcal{R} = (\{1, 2, 3, 4\}, \{a\}, s_{-1}, \text{wt}, (\text{inp}_s, \text{inp}_a), E)$ and $z \in \mathbb{Q}$

$$\text{wt} = \begin{pmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad s_{-1} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1 \end{pmatrix} \quad \text{inp}_s = \text{inp}_a = \begin{pmatrix} -99 \\ -100 \\ 1 \\ 1 \end{pmatrix}$$

States for input $w = a^k$

$$s_t^w = \sigma \left\langle \begin{array}{c} t - 99 \\ t - 100 \\ t + 1 \\ 1 \end{array} \right\rangle \quad \text{for all } 0 \leq t \leq k$$

Example (cont'd)

- States

$$s_t^w = \sigma \left\langle \begin{array}{c} t - 99 \\ t - 100 \\ t + 1 \\ 1 \end{array} \right\rangle$$

for all $0 \leq t \leq k$

- Prediction matrix $E = \begin{pmatrix} z & -z & 0 & -z \\ -z & z & 0 & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

Example (cont'd)

- States

$$s_t^w = \sigma \left\langle \begin{array}{c} t-99 \\ t-100 \\ t+1 \\ 1 \end{array} \right\rangle$$

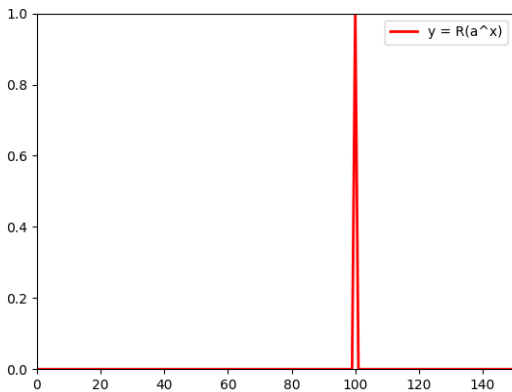
for all $0 \leq t \leq k$

- Prediction matrix $E = \begin{pmatrix} z & -z & 0 & -z \\ -z & z & 0 & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- $\text{out}_t^w = \text{softmax} \left\langle \begin{array}{c} -z \\ 0 \end{array} \right\rangle$ for $0 \leq t \leq 99$

$$\text{out}_t^w = \text{softmax} \left\langle \begin{array}{c} 0 \\ -z \end{array} \right\rangle \text{ for all } 100 \leq t \leq k$$

Example (cont'd)



(graph for $z = 100$)

Theorem [Siegelmann & Sontag 1995]

For every deterministic Turing machine \mathcal{M} there exists RNN

$$(N, \{a\}, s_{-1}, wt, (\text{inp}_\$, \text{inp}_a), E)$$

and 2 designated neurons $n, n' \in N$ such that for all $w = a^k$ and $0 \leq t \leq k$

- $s_t^w(n) - s_t^w(n') \in \{0, 1\}$
- $s_t^w(n) - s_t^w(n') = 1$ iff \mathcal{M} halts on ε after exactly t steps

Definition (Best string)

Best string for RNN $\mathcal{R} = (N, \Sigma, s_{-1}, \text{wt}, (\text{inp}_a)_{a \in \Sigma_S}, E)$ is

$$\hat{w} = \arg \max_{w \in \Sigma^*} \mathcal{R}(w)$$

Best string problem: Is $\mathcal{R}(\hat{w}) \geq \theta$ given \mathcal{R} and $\theta \in [0, 1] \cap \mathbb{Q}$

Definition (Best string)

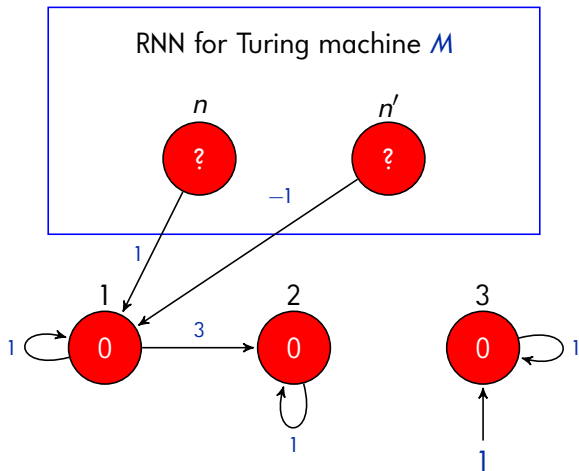
Best string for RNN $\mathcal{R} = (N, \Sigma, s_{-1}, \text{wt}, (\text{inp}_a)_{a \in \Sigma_S}, E)$ is

$$\hat{w} = \arg \max_{w \in \Sigma^*} \mathcal{R}(w)$$

Best string problem: Is $\mathcal{R}(\hat{w}) \geq \theta$ given \mathcal{R} and $\theta \in [0, 1] \cap \mathbb{Q}$

Theorem [CGMMK 2018]

Best string problem undecidable for RNN



Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 1: TM \mathcal{M} does not halt

$$s_t^w = \begin{pmatrix} 0 \\ 0 \\ t+1 \\ \vdots \end{pmatrix}$$

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 1: TM \mathcal{M} does not halt

$$s_t^w = \begin{pmatrix} 0 \\ 0 \\ t+1 \\ \vdots \end{pmatrix}$$

$$\text{out}_t^w(\$) = \frac{e^{-(t+1)}}{e^{-(t+1)} + e^{t+1}} = \frac{1}{1 + e^{2(t+1)}} \leq \frac{1}{1 + e^2} < 0.12$$

Thus $\mathcal{R}(a^k) < 0.12$ for all $k \in \mathbb{N}$

Decidability

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 2: TM M halts after exactly ℓ steps

$$s_t^w = \begin{pmatrix} t \stackrel{?}{>} \ell \\ \sigma \langle 3(t - \ell - 1) \rangle \\ t + 1 \\ \vdots \end{pmatrix}$$

Decidability

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 2: TM M halts after exactly ℓ steps

$$s_t^w = \begin{pmatrix} \overset{?}{t > \ell} \\ \sigma\langle 3(t - \ell - 1) \rangle \\ t + 1 \\ \vdots \end{pmatrix}$$

$$\text{out}_t^w(\$) = \begin{cases} \frac{1}{1+e^{2(t+1)}} & \text{if } t \leq \ell \\ \frac{e^{t-3\ell-5}}{1+e^{t-3\ell-5}} & \text{otherwise} \end{cases}$$

Decidability

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 2: TM M halts after exactly ℓ steps

$$s_t^w = \begin{pmatrix} \overset{?}{t > \ell} \\ \sigma \langle 3(t - \ell - 1) \rangle \\ t + 1 \\ \vdots \end{pmatrix}$$

$$\text{out}_i^w(\$) = \begin{cases} \frac{1}{1+e^{2(t+1)}} & \text{if } t \leq \ell \\ \frac{e^{t-3\ell-5}}{1+e^{t-3\ell-5}} & \text{otherwise} \end{cases}$$

$$\mathcal{R}(a^{3\ell-5}) = \left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) \cdot \left(\prod_{t=\ell+1}^{3\ell-6} \frac{1}{1+e^{t-3\ell-5}} \right) \cdot \frac{1}{2}$$

Definition (Infinite Pochhammer symbol)

Let $x \in \mathbb{R}$ and $q \in [0, 1]$. Infinite q -Pochhammer symbol is

$$(x; q)_{\infty} = \prod_{i \in \mathbb{N}} (1 - xq^i)$$

Definition (Infinite Pochhammer symbol)

Let $x \in \mathbb{R}$ and $q \in [0, 1]$. Infinite q -Pochhammer symbol is

$$(x; q)_{\infty} = \prod_{i \in \mathbb{N}} (1 - xq^i)$$

Examples

$$(-1; e^{-2})_{\infty} = \prod_{i \in \mathbb{N}} (1 + e^{-2i}) \approx 2.32$$

$$(-1; e^{-1})_{\infty} = \prod_{i \in \mathbb{N}} (1 + e^{-i}) \approx 3.36$$

$$\mathcal{R}(a^{3\ell-5}) = \left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1 + e^{2(t+1)}} \right) \cdot \left(\prod_{t=\ell+1}^{3\ell-6} \frac{1}{1 + e^{t-3\ell-5}} \right) \cdot \frac{1}{2}$$

$$\begin{aligned}
 \mathcal{R}(a^{3\ell-5}) &= \left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) \cdot \left(\prod_{t=\ell+1}^{3\ell-6} \frac{1}{1+e^{t-3\ell-5}} \right) \cdot \frac{1}{2} \\
 &= \left(\frac{2}{2} \cdot \prod_{t=1}^{\ell+1} \frac{1}{1+e^{-2t}} \right) \cdot \left(\prod_{t=11}^{2\ell+4} \frac{1}{1+e^{-t}} \right) \cdot \frac{1}{2}
 \end{aligned}$$

$$\begin{aligned}
\mathcal{R}(a^{3\ell-5}) &= \left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) \cdot \left(\prod_{t=\ell+1}^{3\ell-6} \frac{1}{1+e^{t-3\ell-5}} \right) \cdot \frac{1}{2} \\
&= \left(\frac{2}{2} \cdot \prod_{t=1}^{\ell+1} \frac{1}{1+e^{-2t}} \right) \cdot \left(\prod_{t=11}^{2\ell+4} \frac{1}{1+e^{-t}} \right) \cdot \frac{1}{2} \\
&\geq \left(\prod_{t=0}^{\ell+1} \frac{1}{1+e^{-2t}} \right) \cdot 2 \cdot \left(\prod_{t=0}^{2\ell+4} \frac{1}{1+e^{-t}} \right)
\end{aligned}$$

$$\begin{aligned}
\mathcal{R}(a^{3\ell-5}) &= \left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) \cdot \left(\prod_{t=\ell+1}^{3\ell-6} \frac{1}{1+e^{t-3\ell-5}} \right) \cdot \frac{1}{2} \\
&= \left(\frac{2}{2} \cdot \prod_{t=1}^{\ell+1} \frac{1}{1+e^{-2t}} \right) \cdot \left(\prod_{t=11}^{2\ell+4} \frac{1}{1+e^{-t}} \right) \cdot \frac{1}{2} \\
&\geq \left(\prod_{t=0}^{\ell+1} \frac{1}{1+e^{-2t}} \right) \cdot 2 \cdot \left(\prod_{t=0}^{2\ell+4} \frac{1}{1+e^{-t}} \right) \\
&\geq \frac{2}{(-1; e^{-2})_{\infty} \cdot (-1; e^{-1})_{\infty}} \geq 0.25
\end{aligned}$$

So best tree has weight $\hat{w} \geq 0.2$ iff \mathcal{M} halts

Definition (Consistency)

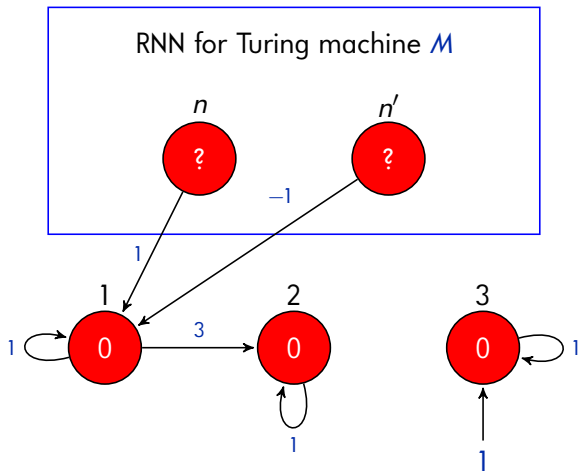
RNN $\mathcal{R} = (N, \Sigma, s_{-1}, wt, (inp_a)_{a \in \Sigma_{\mathcal{S}}}, E)$ is **consistent** if $\sum_{w \in \Sigma^*} \mathcal{R}(w) = 1$

Definition (Consistency)

RNN $\mathcal{R} = (N, \Sigma, s_{-1}, wt, (inp_a)_{a \in \Sigma_S}, E)$ is **consistent** if $\sum_{w \in \Sigma^*} \mathcal{R}(w) = 1$

Theorem [CGMMK 2018]

Consistency undecidable for RNN



Decidability

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 1: TM M does not halt

$$s_t^w = (0 \quad 0 \quad t+1 \quad \dots)^T$$

$$\text{out}_t^w(\$) = \frac{1}{1 + e^{2(t+1)}}$$

Decidability

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 1: TM M does not halt

$$s_t^w = (0 \quad 0 \quad t+1 \quad \cdots)^T \qquad \text{out}_t^w(\$) = \frac{1}{1 + e^{2(t+1)}}$$

Thus \mathcal{R} is inconsistent because

$$\mathcal{R}(\alpha^k) = \frac{1}{1 + e^{2(k+1)}} \cdot \prod_{t=0}^{k-1} \frac{e^{2(t+1)}}{1 + e^{2(t+1)}}$$

Decidability

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 1: TM M does not halt

$$s_t^w = (0 \quad 0 \quad t+1 \quad \cdots)^T \quad \text{out}_t^w(\$) = \frac{1}{1 + e^{2(t+1)}}$$

Thus \mathcal{R} is inconsistent because

$$\mathcal{R}(\alpha^k) = \frac{1}{1 + e^{2(k+1)}} \cdot \prod_{t=0}^{k-1} \frac{e^{2(t+1)}}{1 + e^{2(t+1)}}$$
$$\sum_{k \in \mathbb{N}} \mathcal{R}(\alpha^k) = 1 - \prod_{t \in \mathbb{N}} \frac{1}{1 + e^{-2(t+1)}} = 1 - \frac{2}{(-1; e^{-2})_\infty} < 0.14$$

Decidability

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 2: TM M halts after exactly l steps

$$s_t^w = \begin{pmatrix} \overset{?}{t > l} \\ \sigma\langle 3(t-l-1) \rangle \\ t+1 \\ \vdots \end{pmatrix} \quad \text{out}_t^w(\$) = \begin{cases} \frac{1}{1+e^{2(t+1)}} & \text{if } t \leq l \\ \frac{e^{t-3l-5}}{1+e^{t-3l-5}} & \text{otherwise} \end{cases}$$

Prediction matrix $E = \begin{pmatrix} 0 & 1 & -1 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \end{pmatrix}$ $\begin{matrix} \$ \\ a \end{matrix}$

- Case 2: TM M halts after exactly ℓ steps

$$s_t^w = \begin{pmatrix} \overset{?}{t > \ell} \\ \sigma\langle 3(t - \ell - 1) \rangle \\ t + 1 \\ \vdots \end{pmatrix} \quad \text{out}_t^w(\$) = \begin{cases} \frac{1}{1 + e^{2(t+1)}} & \text{if } t \leq \ell \\ \frac{e^{t-3\ell-5}}{1 + e^{t-3\ell-5}} & \text{otherwise} \end{cases}$$

Hence \mathcal{R} is consistent because

$$\sum_{k \in \mathbb{N}} \mathcal{R}(a^k) = 1$$

(Appendix Lemma B)

Notes

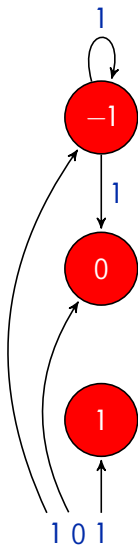
- Consistency undecidable by previous theorem
- Backpropagation-through-time (BPTT) might return inconsistent RNN

	Best path	Best string
RNN	Undecidable	
Consistent RNN	NP-c[†]	
Det. WA/PCFG	P	
WA/PCFG	NP-c	

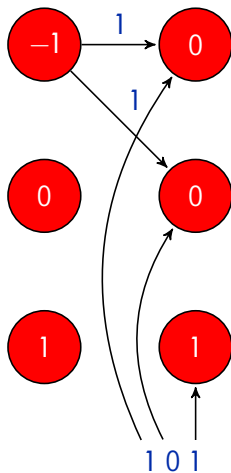
[†] restricted to solutions of polynomial length

Training

Unrolling

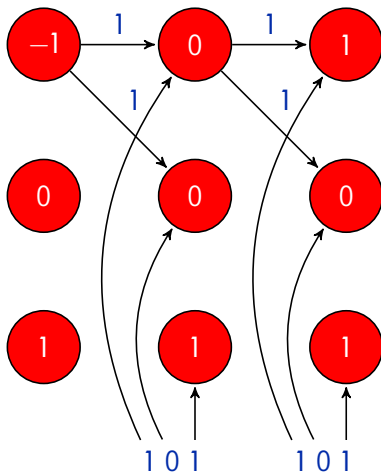


Unrolling



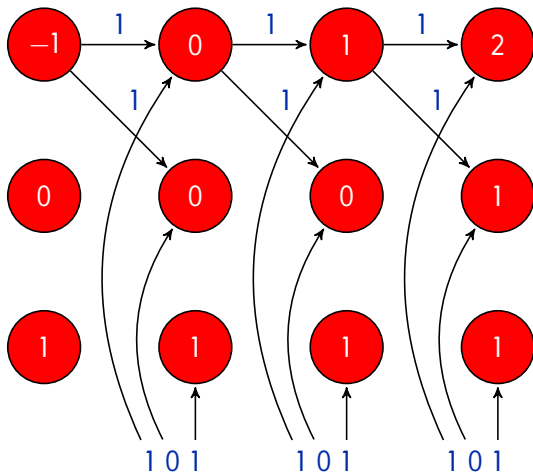
Training

Unrolling



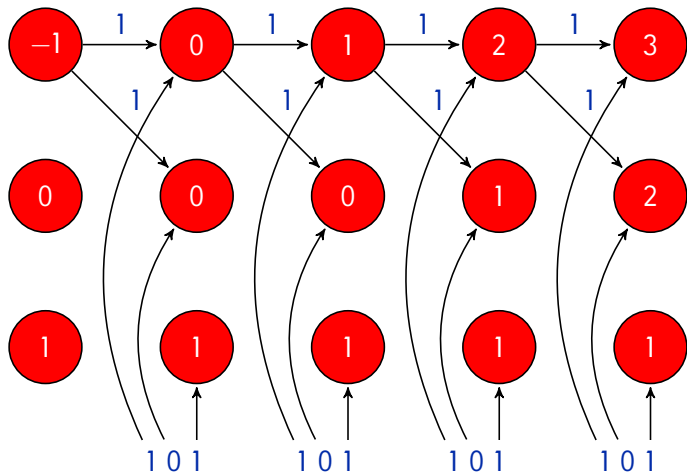
Training

Unrolling



Training

Unrolling



Loss function

- Measures errors currently made
- **Cross-entropy loss** for predictions p

Loss function

- Measures errors currently made
- **Cross-entropy loss** for predictions p
 - ▶ Next letter given input w with reference a

$$L(w, a, p) = -\log p(w, a)$$

Loss function

- Measures errors currently made
- **Cross-entropy loss** for predictions p
 - ▶ Next letter given input w with reference a

$$L(w, a, p) = -\log p(w, a)$$

- ▶ Input $w \in \mathcal{C}$

$$L(w, p) = \text{Avg}_{1 \leq i \leq |w|} L(w[1] \cdots w[i-1], w[i], p)$$

Loss function

- Measures errors currently made
- **Cross-entropy loss** for predictions p
 - ▶ Next letter given input w with reference a

$$L(w, a, p) = -\log p(w, a)$$

- ▶ Input $w \in \mathcal{C}$

$$L(w, p) = \text{Avg}_{1 \leq i \leq |w|} L(w[1] \cdots w[i-1], w[i], p)$$

- ▶ Language \mathcal{C}

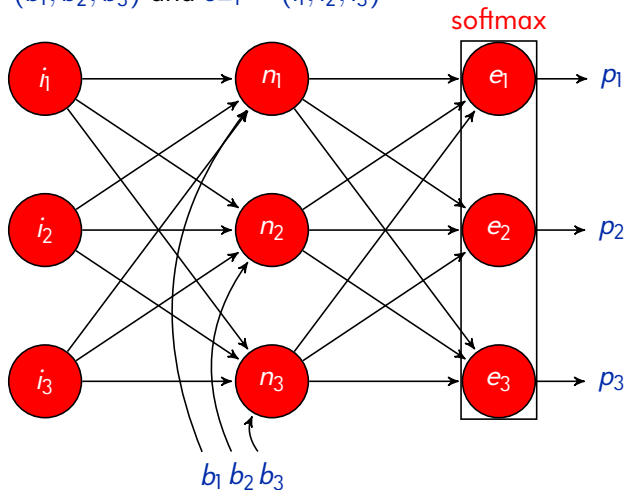
$$L(\mathcal{C}, p) = \text{Avg}_{w \in \mathcal{C}} L(w, p)$$

Feed-forward network

- First letter, $|\Sigma| = 2$ letters ($|\Sigma_{\mathcal{S}}| = 3$ outputs), $|N| = 3$ neurons
- $\text{inp}_{\mathcal{S}} = (b_1, b_2, b_3)$ and $s_{-1} = (i_1, i_2, i_3)$

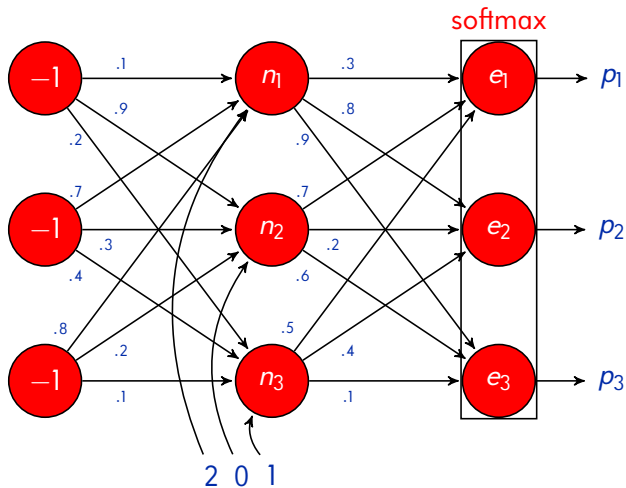
Feed-forward network

- First letter, $|\Sigma| = 2$ letters ($|\Sigma_{\S}| = 3$ outputs), $|N| = 3$ neurons
- $\text{inp}_{\S} = (b_1, b_2, b_3)$ and $s_{-1} = (i_1, i_2, i_3)$



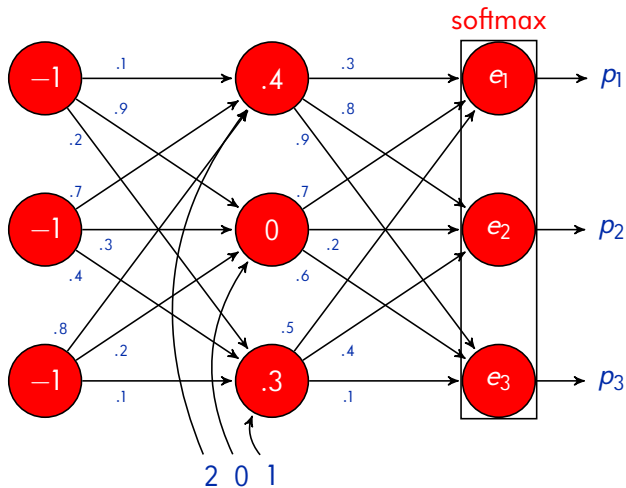
Training

Forward pass



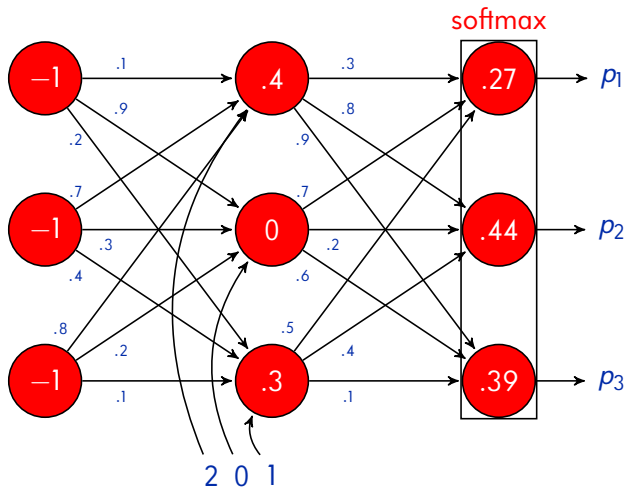
Training

Forward pass



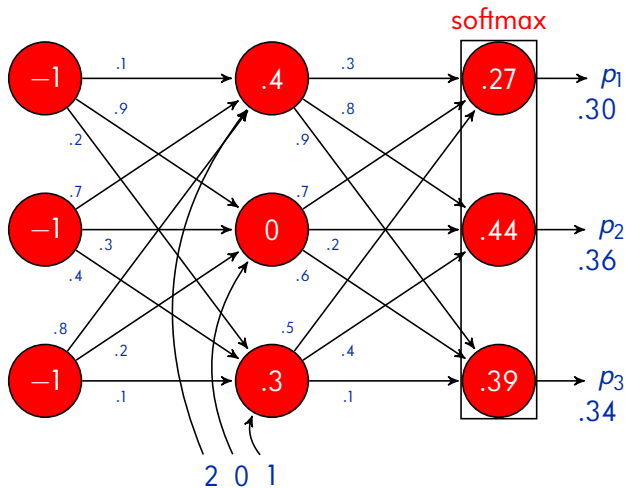
Training

Forward pass



Training

Forward pass



Calculation of loss

- Obtain output probabilities $(p_1, p_2, p_3) \approx (.30, .36, .34)$
- Loss $L = -\log(p_3) \approx 1.08$ for correct output corresponding to p_3

Calculation of loss

- Obtain output probabilities $(p_1, p_2, p_3) \approx (.30, .36, .34)$
- Loss $L = -\log(p_3) \approx 1.08$ for correct output corresponding to p_3

Assigning blame and computing changes

- Want lower loss, so calculate partial derivatives of loss function

Calculation of loss

- Obtain output probabilities $(p_1, p_2, p_3) \approx (.30, .36, .34)$
- Loss $L = -\log(p_3) \approx 1.08$ for correct output corresponding to p_3

Assigning blame and computing changes

- Want lower loss, so calculate partial derivatives of loss function
- Treat all edge weights $wt(\cdot, \cdot)$ as variables (here: entries of wt and E)
- Calculate $\frac{\partial L}{\partial wt(n_1, e_3)}$ with $L = -\log p_3$

Computing derivatives

$$L = -\log p_3 \quad p_3 = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}$$

$$e_3 = \max(0, z_3) \quad z_3 = n_1 \cdot \text{wt}(n_1, e_3) + n_2 \cdot \text{wt}(n_2, e_3) + n_3 \cdot \text{wt}(n_3, e_3)$$

Computing derivatives

$$L = -\log p_3 \quad p_3 = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}$$

$$e_3 = \max(0, z_3) \quad z_3 = n_1 \cdot \text{wt}(n_1, e_3) + n_2 \cdot \text{wt}(n_2, e_3) + n_3 \cdot \text{wt}(n_3, e_3)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{\partial L}{\partial p_3} \frac{\partial p_3}{\partial e_3} \frac{\partial e_3}{\partial z_3} \frac{\partial z_3}{\partial \text{wt}(n_1, e_3)}$$

Computing derivatives

$$L = -\log p_3 \quad p_3 = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}$$

$$e_3 = \max(0, z_3) \quad z_3 = n_1 \cdot \text{wt}(n_1, e_3) + n_2 \cdot \text{wt}(n_2, e_3) + n_3 \cdot \text{wt}(n_3, e_3)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{\partial L}{\partial p_3} \frac{\partial p_3}{\partial e_3} \frac{\partial e_3}{\partial z_3} \frac{\partial z_3}{\partial \text{wt}(n_1, e_3)}$$

- $\frac{\partial L}{\partial p_3} = -p_3^{-1}$

Computing derivatives

$$L = -\log p_3 \quad p_3 = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}$$

$$e_3 = \max(0, z_3) \quad z_3 = n_1 \cdot \text{wt}(n_1, e_3) + n_2 \cdot \text{wt}(n_2, e_3) + n_3 \cdot \text{wt}(n_3, e_3)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{\partial L}{\partial p_3} \frac{\partial p_3}{\partial e_3} \frac{\partial e_3}{\partial z_3} \frac{\partial z_3}{\partial \text{wt}(n_1, e_3)}$$

- $\frac{\partial L}{\partial p_3} = -p_3^{-1}$
- $\frac{\partial p_3}{\partial e_3} = \frac{e^{e_3}(e^{e_1} + e^{e_2} + e^{e_3}) - e^{e_3} \cdot e^{e_3}}{(e^{e_1} + e^{e_2} + e^{e_3})^2} = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}} - \left(\frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}\right)^2 = p_3(1 - p_3)$

Computing derivatives

$$L = -\log p_3 \quad p_3 = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}$$

$$e_3 = \max(0, z_3) \quad z_3 = n_1 \cdot \text{wt}(n_1, e_3) + n_2 \cdot \text{wt}(n_2, e_3) + n_3 \cdot \text{wt}(n_3, e_3)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{\partial L}{\partial p_3} \frac{\partial p_3}{\partial e_3} \frac{\partial e_3}{\partial z_3} \frac{\partial z_3}{\partial \text{wt}(n_1, e_3)}$$

- $\frac{\partial L}{\partial p_3} = -p_3^{-1}$
- $\frac{\partial p_3}{\partial e_3} = \frac{e^{e_3}(e^{e_1} + e^{e_2} + e^{e_3}) - e^{e_3} \cdot e^{e_3}}{(e^{e_1} + e^{e_2} + e^{e_3})^2} = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}} - \left(\frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}\right)^2 = p_3(1 - p_3)$
- $\frac{\partial e_3}{\partial z_3} = 1$ if $z_3 > 0$ and 0 otherwise (well almost)

Computing derivatives

$$L = -\log p_3 \quad p_3 = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}$$

$$e_3 = \max(0, z_3) \quad z_3 = n_1 \cdot \text{wt}(n_1, e_3) + n_2 \cdot \text{wt}(n_2, e_3) + n_3 \cdot \text{wt}(n_3, e_3)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{\partial L}{\partial p_3} \frac{\partial p_3}{\partial e_3} \frac{\partial e_3}{\partial z_3} \frac{\partial z_3}{\partial \text{wt}(n_1, e_3)}$$

- $\frac{\partial L}{\partial p_3} = -p_3^{-1}$

- $\frac{\partial p_3}{\partial e_3} = \frac{e^{e_3}(e^{e_1} + e^{e_2} + e^{e_3}) - e^{e_3} \cdot e^{e_3}}{(e^{e_1} + e^{e_2} + e^{e_3})^2} = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}} - \left(\frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}\right)^2 = p_3(1 - p_3)$

- $\frac{\partial e_3}{\partial z_3} = 1$ if $z_3 > 0$ and 0 otherwise (well almost)

- $\frac{\partial z_3}{\partial \text{wt}(n_1, e_3)} = n_1$

Computing derivatives

$$L = -\log p_3 \quad p_3 = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}$$

$$e_3 = \max(0, z_3) \quad z_3 = n_1 \cdot \text{wt}(n_1, e_3) + n_2 \cdot \text{wt}(n_2, e_3) + n_3 \cdot \text{wt}(n_3, e_3)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{\partial L}{\partial p_3} \frac{\partial p_3}{\partial e_3} \frac{\partial e_3}{\partial z_3} \frac{\partial z_3}{\partial \text{wt}(n_1, e_3)} = -\frac{p_3}{p_3} (1 - p_3) n_1 = n_1 (p_3 - 1)$$

- $\frac{\partial L}{\partial p_3} = -p_3^{-1}$

- $\frac{\partial p_3}{\partial e_3} = \frac{e^{e_3}(e^{e_1} + e^{e_2} + e^{e_3}) - e^{e_3} \cdot e^{e_3}}{(e^{e_1} + e^{e_2} + e^{e_3})^2} = \frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}} - \left(\frac{e^{e_3}}{e^{e_1} + e^{e_2} + e^{e_3}}\right)^2 = p_3(1 - p_3)$

- $\frac{\partial e_3}{\partial z_3} = 1$ if $z_3 > 0$ and 0 otherwise (well almost)

- $\frac{\partial z_3}{\partial \text{wt}(n_1, e_3)} = n_1$

Multiple inputs

- Superscript in parentheses indicating input
e.g. $n_3^{(i)}$ for n_3 of i -th input
- Indicator bit $f_k^{(i)} \in \{0, 1\}$ indicating
whether desired output symbol for i -th input corresponds to p_k

Computing derivatives for ℓ inputs

$$L = -\frac{1}{\ell} \sum_{i=1}^{\ell} \sum_{k=1}^3 t_k^{(i)} \log p_k^{(i)}$$

$$e_m^{(i)} = \max(0, z_m^{(i)})$$

$$p_k^{(i)} = \frac{e^{e_k^{(i)}}}{e^{e_1^{(i)}} + e^{e_2^{(i)}} + e^{e_3^{(i)}}}$$

$$z_m^{(i)} = \sum_{h=1}^3 n_h^{(i)} \cdot \text{wt}(n_h, e_m)$$

Computing derivatives for ℓ inputs

$$L = -\frac{1}{\ell} \sum_{i=1}^{\ell} \sum_{k=1}^3 t_k^{(i)} \log p_k^{(i)}$$

$$p_k^{(i)} = \frac{e^{e_k^{(i)}}}{e^{e_1^{(i)}} + e^{e_2^{(i)}} + e^{e_3^{(i)}}}$$

$$e_m^{(i)} = \max(0, z_m^{(i)})$$

$$z_m^{(i)} = \sum_{h=1}^3 n_h^{(i)} \cdot \text{wt}(n_h, e_m)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \sum_{i=1}^{\ell} \sum_{k=1}^3 \sum_{m=1}^3 \frac{\partial L}{\partial p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial e_m^{(i)}} \frac{\partial e_m^{(i)}}{\partial z_m^{(i)}} \frac{\partial z_m^{(i)}}{\partial \text{wt}(n_1, e_3)}$$

Computing derivatives for ℓ inputs

$$L = -\frac{1}{\ell} \sum_{i=1}^{\ell} \sum_{k=1}^3 t_k^{(i)} \log p_k^{(i)}$$

$$p_k^{(i)} = \frac{e^{e_k^{(i)}}}{e^{e_1^{(i)}} + e^{e_2^{(i)}} + e^{e_3^{(i)}}}$$

$$e_m^{(i)} = \max(0, z_m^{(i)})$$

$$z_m^{(i)} = \sum_{h=1}^3 n_h^{(i)} \cdot \text{wt}(n_h, e_m)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \sum_{i=1}^{\ell} \sum_{k=1}^3 \sum_{m=1}^3 \frac{\partial L}{\partial p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial e_m^{(i)}} \frac{\partial e_m^{(i)}}{\partial z_m^{(i)}} \frac{\partial z_m^{(i)}}{\partial \text{wt}(n_1, e_3)}$$

- $\frac{\partial L}{\partial p_k^{(i)}} = -\frac{t_k^{(i)}}{\ell p_k^{(i)}}$

Computing derivatives for ℓ inputs

$$L = -\frac{1}{\ell} \sum_{i=1}^{\ell} \sum_{k=1}^3 t_k^{(i)} \log p_k^{(i)}$$

$$p_k^{(i)} = \frac{e^{e_k^{(i)}}}{e^{e_1^{(i)}} + e^{e_2^{(i)}} + e^{e_3^{(i)}}}$$

$$e_m^{(i)} = \max(0, z_m^{(i)})$$

$$z_m^{(i)} = \sum_{h=1}^3 n_h^{(i)} \cdot \text{wt}(n_h, e_m)$$

- Chain rule

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \sum_{i=1}^{\ell} \sum_{k=1}^3 \sum_{m=1}^3 \frac{\partial L}{\partial p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial e_m^{(i)}} \frac{\partial e_m^{(i)}}{\partial z_m^{(i)}} \frac{\partial z_m^{(i)}}{\partial \text{wt}(n_1, e_3)}$$

- $\frac{\partial L}{\partial p_k^{(i)}} = -\frac{t_k^{(i)}}{\ell p_k^{(i)}}$
- $\frac{\partial e_m^{(i)}}{\partial z_m^{(i)}} = 1$ if $z_m^{(i)} > 0$ and 0 otherwise

(well almost)

Computing derivatives for ℓ inputs

$$p_k^{(i)} = \frac{e^{e_k^{(i)}}}{e^{e_1^{(i)}} + e^{e_2^{(i)}} + e^{e_3^{(i)}}}$$

$$z_m^{(i)} = \sum_{h=1}^3 n_h^{(i)} \cdot \text{wt}(n_h, e_m)$$

$$\frac{\partial p_k^{(i)}}{\partial e_m^{(i)}} = \begin{cases} p_k^{(i)}(1 - p_k^{(i)}) & \text{if } k = m \\ -p_k^{(i)} p_m^{(i)} & \text{otherwise} \end{cases}$$

$$\frac{\partial z_m^{(i)}}{\partial \text{wt}(n_1, e_3)} = \begin{cases} n_1^{(i)} & \text{if } m = 3 \\ 0 & \text{otherwise} \end{cases}$$

Computing derivatives for ℓ inputs

$$p_k^{(i)} = \frac{e^{e_k^{(i)}}}{e^{e_1^{(i)}} + e^{e_2^{(i)}} + e^{e_3^{(i)}}} \quad z_m^{(i)} = \sum_{h=1}^3 n_h^{(i)} \cdot \text{wt}(n_h, e_m)$$

$$\frac{\partial p_k^{(i)}}{\partial e_m^{(i)}} = \begin{cases} p_k^{(i)}(1 - p_k^{(i)}) & \text{if } k = m \\ -p_k^{(i)} p_m^{(i)} & \text{otherwise} \end{cases}$$

$$\frac{\partial z_m^{(i)}}{\partial \text{wt}(n_1, e_3)} = \begin{cases} n_1^{(i)} & \text{if } m = 3 \\ 0 & \text{otherwise} \end{cases}$$

Thus

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{1}{\ell} \sum_{i \in \{1, \dots, \ell\}, z_3^{(i)} > 0} (p_3^{(i)} - t_3^{(i)}) n_1^{(i)} \quad (\text{Appendix Lemma C})$$

Gradient descent

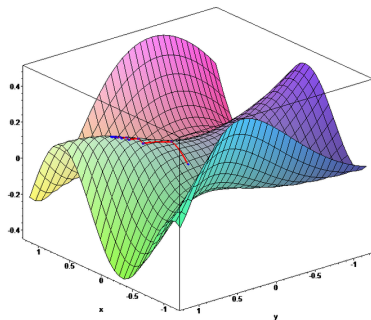
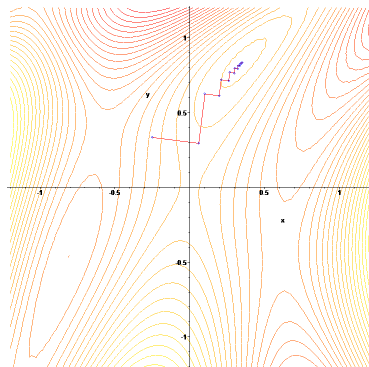
- $f: \mathbb{R}^n \rightarrow \mathbb{R}$ defined & differentiable around $a \in \mathbb{R}^n$ decreases fastest going from a in direction of gradient $-\nabla f(a)$

Gradient descent

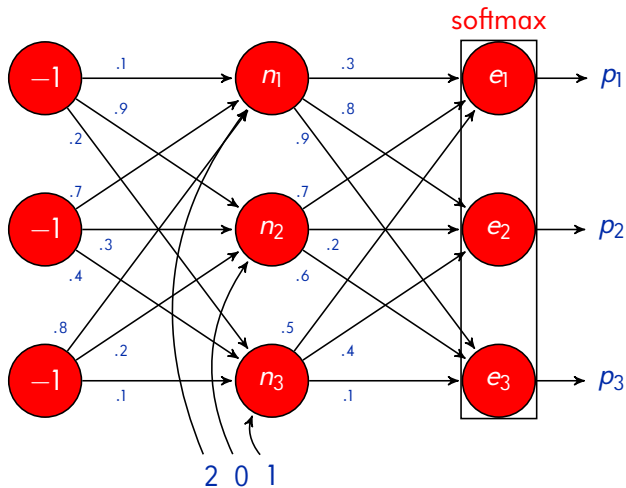
- $f: \mathbb{R}^n \rightarrow \mathbb{R}$ defined & differentiable around $a \in \mathbb{R}^n$ decreases fastest going from a in direction of gradient $-\nabla f(a)$
- $a_{n+1} = a_n - \gamma \nabla f(a_n)$ with $\gamma \in \mathbb{R}_+$ small yields $f(a_{n+1}) \leq f(a_n)$

Gradient descent

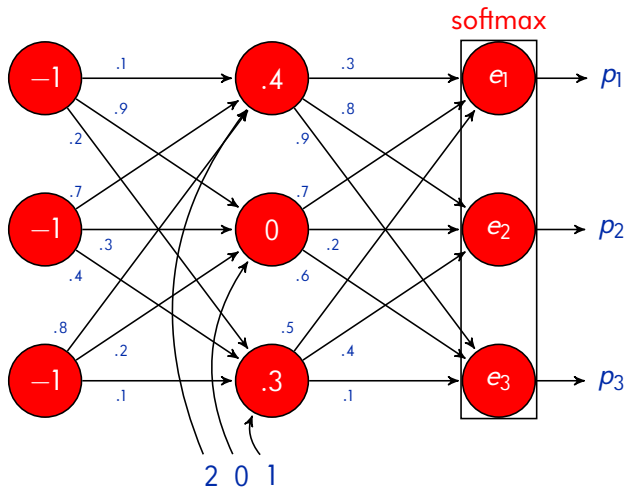
- $f: \mathbb{R}^n \rightarrow \mathbb{R}$ defined & differentiable around $a \in \mathbb{R}^n$ decreases fastest going from a in direction of gradient $-\nabla f(a)$
- $a_{n+1} = a_n - \gamma \nabla f(a_n)$ with $\gamma \in \mathbb{R}_+$ small yields $f(a_{n+1}) \leq f(a_n)$
- **Learning rate** γ (often varied over time)



Backpropagation

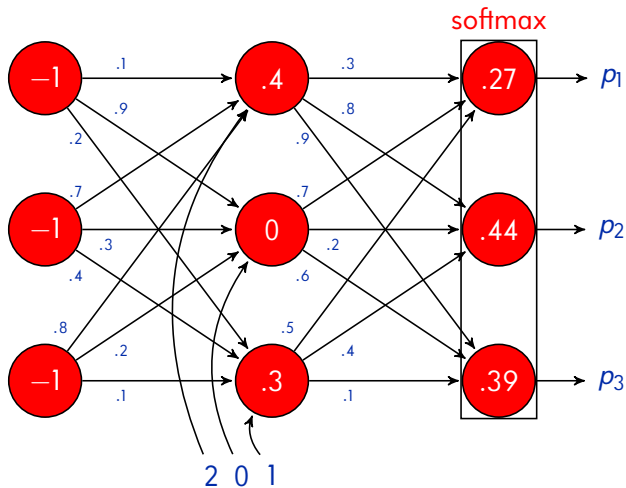


Backpropagation



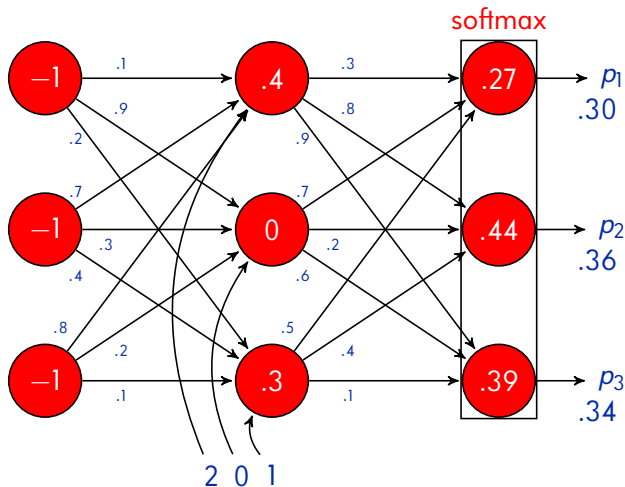
Training

Backpropagation



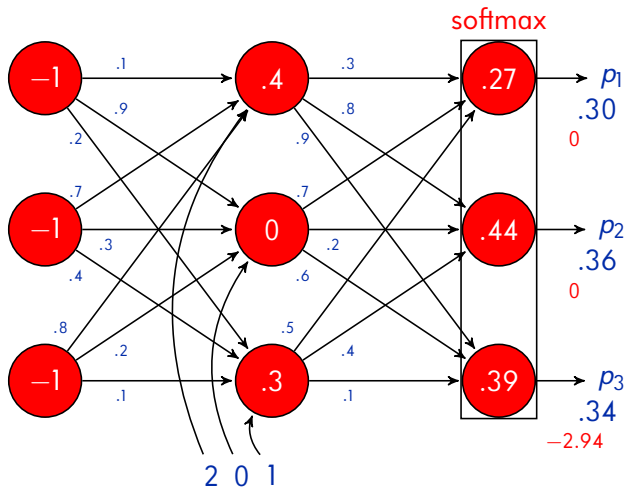
Backpropagation

- $\frac{\partial L}{\partial p_1} = \frac{\partial L}{\partial p_2} = 0$ and $\frac{\partial L}{\partial p_3} = -p_3^{-1} \approx -2.94$



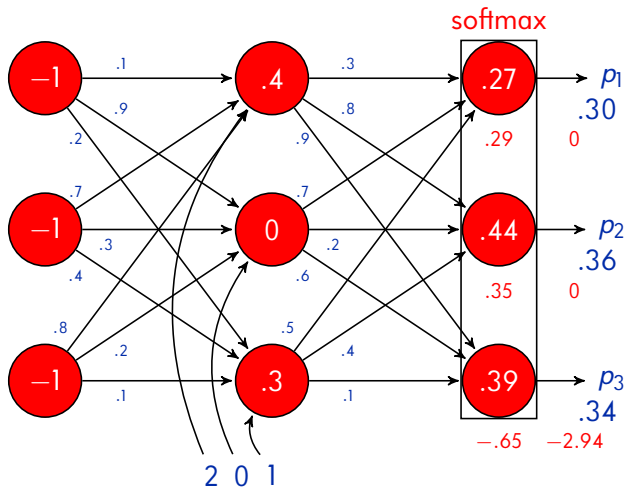
Backpropagation

- $\frac{\partial p_3}{\partial e_1} = -p_3 p_1 \approx -0.10$ and $\frac{\partial p_3}{\partial e_3} = p_3(1 - p_3) \approx 0.22$



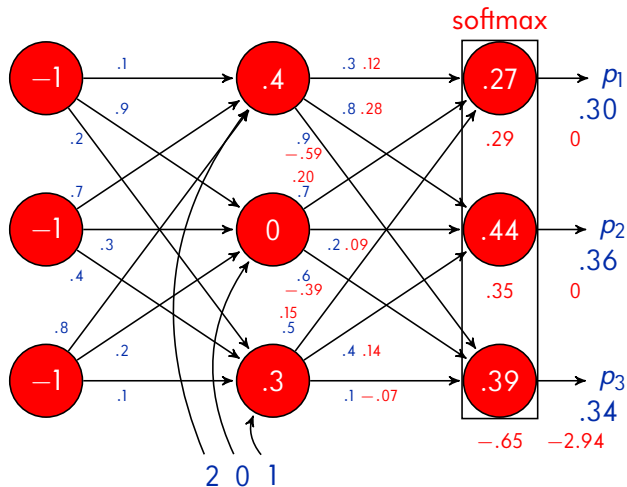
Backpropagation

- $\frac{\partial z_m}{\partial \text{wt}(n_i, e_m)} = n_i$ for all $i, m \in \{1, 2, 3\}$



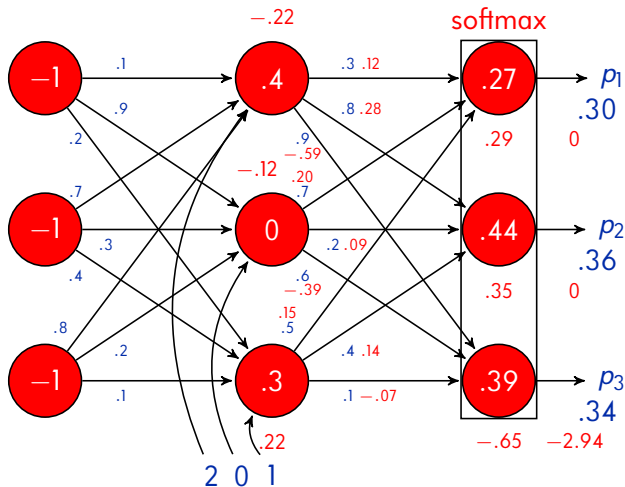
Backpropagation

- $\frac{\partial z_m}{\partial n_i} = wt(n_i, e_m)$ for all $i, m \in \{1, 2, 3\}$



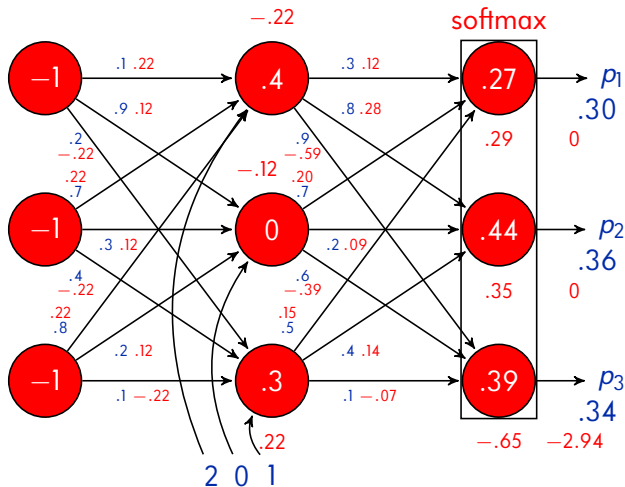
Backpropagation

- $\frac{\partial n_j}{\partial \text{wt}(i_k, n_j)} = i_k$ for all $i, k \in \{1, 2, 3\}$



Backpropagation

- $\frac{\partial n_j}{\partial \text{wt}(i_k, n_j)} = i_k$ for all $i, k \in \{1, 2, 3\}$



Partial derivatives

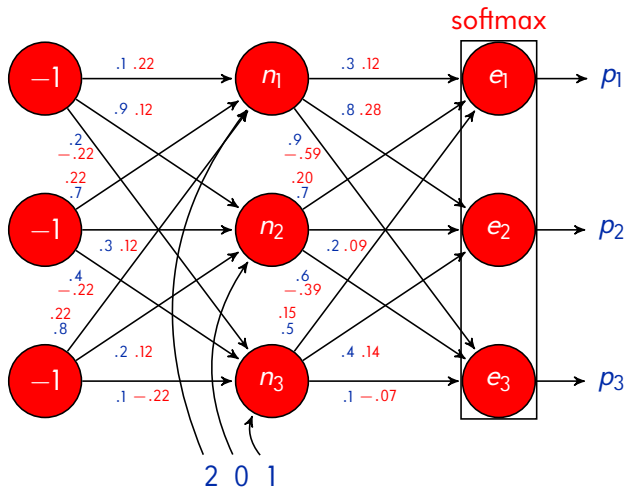
$$\nabla L = \begin{pmatrix} \partial L / \partial \text{wt}(i_1, n_1) \\ \partial L / \partial \text{wt}(i_1, n_2) \\ \dots \\ \partial L / \partial \text{wt}(i_3, n_3) \\ \partial L / \partial \text{wt}(n_1, e_1) \\ \partial L / \partial \text{wt}(n_1, e_2) \\ \dots \\ \partial L / \partial \text{wt}(n_3, e_3) \end{pmatrix} \approx \begin{pmatrix} .22 \\ .12 \\ \dots \\ -.22 \\ .12 \\ .28 \\ \dots \\ -.07 \end{pmatrix}$$

- Learning rate $\gamma = 1$
- Compute $\text{wt}' = \text{wt} - \gamma \nabla L$

(unrealistically large)

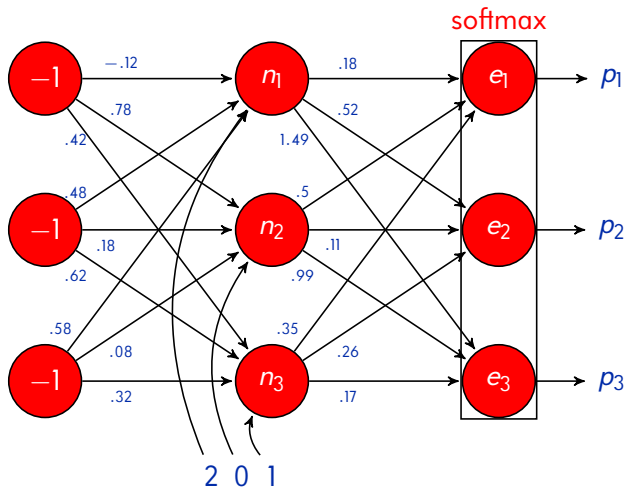
Backpropagation

- Previous loss $-\log(0.34) \approx 1.08$



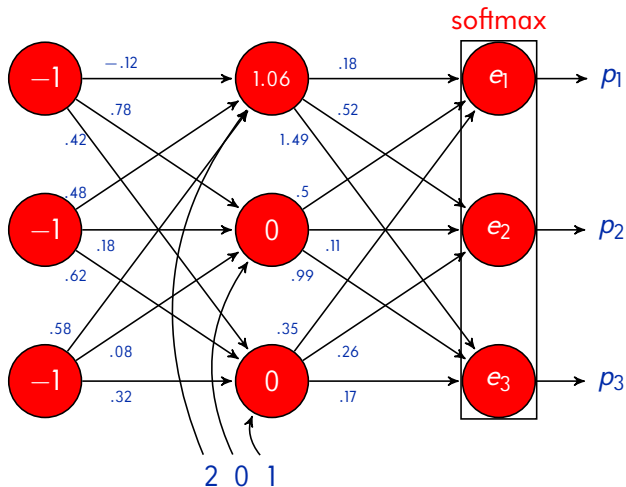
Backpropagation

- Previous loss $-\log(0.34) \approx 1.08$



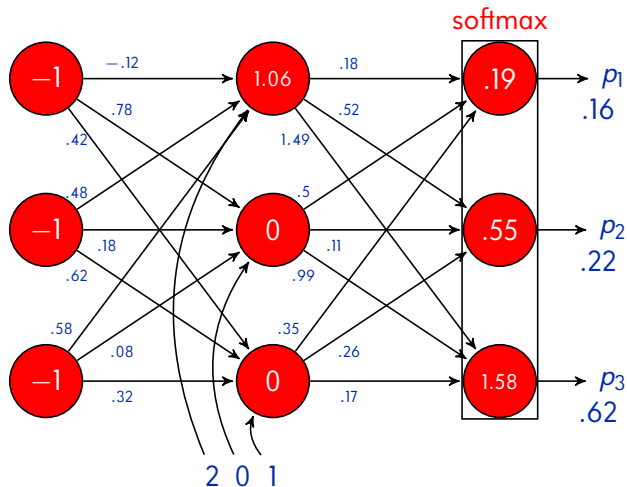
Backpropagation

- Previous loss $-\log(0.34) \approx 1.08$



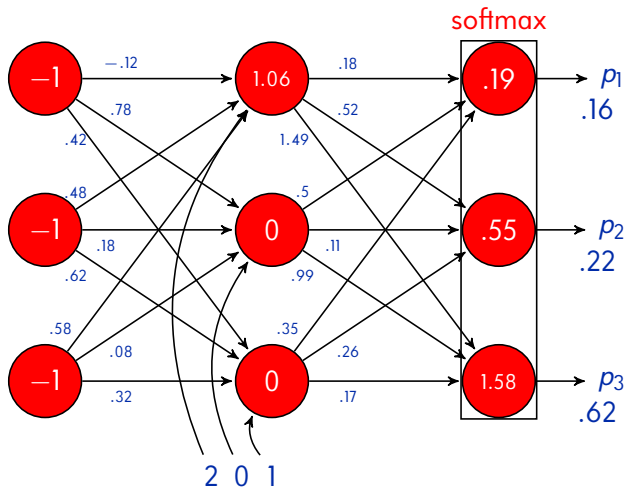
Backpropagation

- Previous loss $-\log(0.34) \approx 1.08$



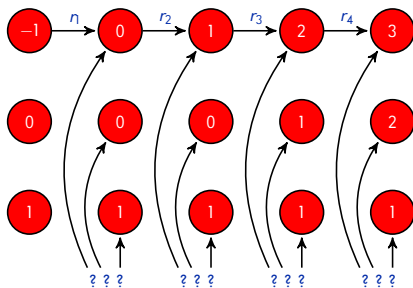
Backpropagation

- Previous loss $-\log(0.34) \approx 1.08$
- New loss $-\log(0.62) \approx 0.48$



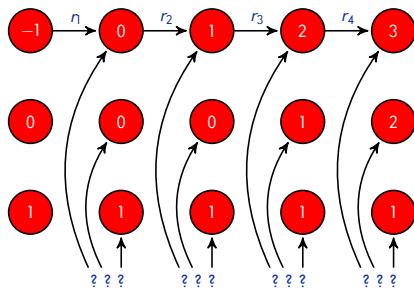
Backpropagation through time

- Train feed-forward network and obtain $\frac{\partial L}{\partial r_i}$ for all $i \in \{1, 2, 3, 4\}$



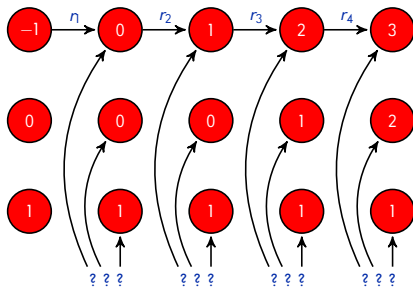
Backpropagation through time

- Train feed-forward network and obtain $\frac{\partial L}{\partial r_i}$ for all $i \in \{1, 2, 3, 4\}$
- $r_1 = r_2 = r_3 = r_4$ due to unrolling from weight r



Backpropagation through time

- Train feed-forward network and obtain $\frac{\partial L}{\partial r_i}$ for all $i \in \{1, 2, 3, 4\}$
- $r_1 = r_2 = r_3 = r_4$ due to unrolling from weight r
- Update $r' = r - \frac{\gamma}{4} \left(\sum_{i=1}^4 \frac{\partial L}{\partial r_i} \right)$ by average of updates for copies



Stochastic gradient descent

- Computing gradients for training data usually too slow & expensive
- Cut data into smaller chunks (randomize order, equal size parts)

Stochastic gradient descent

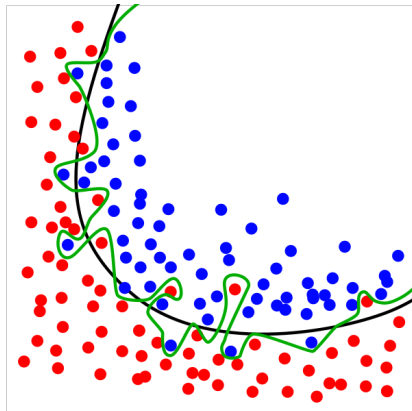
- Computing gradients for training data usually too slow & expensive
- Cut data into smaller chunks (randomize order, equal size parts)
- **Batch** = chunk for which gradients are computed together (e.g. 256 training examples)

Stochastic gradient descent

- Computing gradients for training data usually too slow & expensive
- Cut data into smaller chunks (randomize order, equal size parts)
- **Batch** = chunk for which gradients are computed together (e.g. 256 training examples)
- Train for d **epochs** = d iterations over all training examples

Final notes

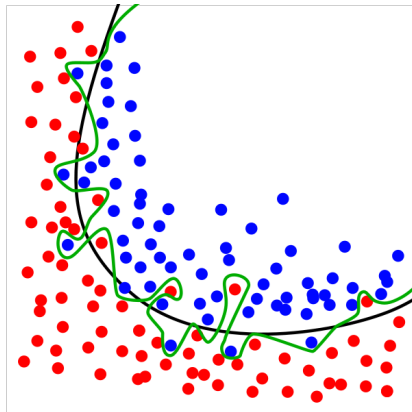
- Too few epochs \rightarrow bad data fit
- Too many epochs \rightarrow overfitting



© Chabacano

Final notes

- Too few epochs → bad data fit
- Too many epochs → overfitting
- Techniques to avoid overfitting: validation loss, drop-out, ...



© Chabacano

That's all folks!

Thank you for the attention

Lemma A

$$\prod_{t \in \mathbb{N}} \frac{e^{k(t+1)}}{1 + e^{k(t+1)}} = \frac{2}{(-1; e^{-k})_{\infty}}$$

for all $k \in \mathbb{N}$

Additional lemma

Lemma A

$$\prod_{t \in \mathbb{N}} \frac{e^{k(t+1)}}{1 + e^{k(t+1)}} = \frac{2}{(-1; e^{-k})_{\infty}} \quad \text{for all } k \in \mathbb{N}$$

Proof

$$\begin{aligned} \prod_{t \in \mathbb{N}} \frac{e^{k(t+1)}}{1 + e^{k(t+1)}} &= \prod_{t \in \mathbb{N}_+} \frac{e^{kt}}{1 + e^{kt}} \cdot \frac{e^{-kt}}{e^{-kt}} = \prod_{t \in \mathbb{N}_+} \frac{\cancel{e^{kt}}}{1 + e^{kt}} \cdot \frac{\cancel{e^{-kt}}}{e^{-kt}} \\ &= \prod_{t \in \mathbb{N}_+} \frac{1}{1 + e^{-kt}} = 2 \cdot \prod_{t \in \mathbb{N}} \frac{1}{1 + e^{-kt}} = \frac{2}{(-1; e^{-k})_{\infty}} \end{aligned}$$

Additional lemma

Lemma B

$$\sum_{k \in \mathbb{N}} \mathcal{R}(a^k) = 1$$

Additional lemma

Lemma B

$$\sum_{k \in \mathbb{N}} \mathcal{R}(a^k) = 1$$

Proof

$$\begin{aligned} \sum_{k \in \mathbb{N}} \mathcal{R}(a^k) &= \left(\sum_{k=0}^{\ell} \mathcal{R}(a^k) \right) + \left(\sum_{k=\ell+1}^{\infty} \mathcal{R}(a^k) \right) \\ &= \sum_{k=0}^{\ell} \left(\frac{1}{1+e^{2(k+1)}} \cdot \prod_{t=0}^{k-1} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) + \sum_{k=\ell+1}^{\infty} \frac{e^{k-3\ell-5}}{1+e^{k-3\ell-5}} \cdot \left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) \cdot \left(\prod_{t=\ell+1}^{k-1} \frac{1}{1+e^{t-3\ell-5}} \right) \\ &= 1 - \left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) + \underbrace{\left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) \cdot \sum_{k=\ell+1}^{\infty} \frac{e^{k-3\ell-5}}{1+e^{k-3\ell-5}} \cdot \left(\prod_{t=\ell+1}^{k-1} \frac{1}{1+e^{t-3\ell-5}} \right)}_{1 - \prod_{t=\ell+1}^{\infty} \frac{1}{1+e^{t-3\ell-5}}} \\ &= 1 - \left(\prod_{t=0}^{\ell} \frac{e^{2(t+1)}}{1+e^{2(t+1)}} \right) \cdot \left(1 - 1 + \prod_{t=\ell+1}^{\infty} \frac{1}{1+e^{t-3\ell-5}} \right) = 1 - \left(\prod_{t=0}^{\ell} \frac{1}{1+e^{-2(t+1)}} \right) \cdot \left(\prod_{t=\ell+1}^{\infty} \frac{1}{1+e^{t-3\ell-5}} \right) \\ &\geq 1 - \left(\prod_{t=0}^{\ell} \frac{1}{1+e^{-2(t+1)}} \right) \cdot \left(\prod_{t \in \mathbb{N}} \frac{1}{1+e^t} \right) = 1 \end{aligned}$$

Lemma C

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{1}{\ell} \sum_{i \in \{1, \dots, \ell\}, z_3^{(i)} > 0} (p_3^{(i)} - t_3^{(i)}) n_1^{(i)}$$

Additional lemma

Lemma C

$$\frac{\partial L}{\partial \text{wt}(n_1, e_3)} = \frac{1}{\ell} \sum_{i \in \{1, \dots, \ell\}, z_3^{(i)} > 0} (p_3^{(i)} - t_3^{(i)}) n_1^{(i)}$$

Proof

$$\begin{aligned} \frac{\partial L}{\partial \text{wt}(n_1, e_3)} &= \sum_{i=1}^{\ell} \sum_{k=1}^3 \sum_{m=1}^3 \frac{\partial L}{\partial p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial e_m^{(i)}} \frac{\partial e_m^{(i)}}{\partial z_m^{(i)}} \frac{\partial z_m^{(i)}}{\partial \text{wt}(n_1, e_3)} = \sum_{i \in \{1, \dots, \ell\}, z_3^{(i)} > 0} \sum_{j=1}^3 -\frac{t_k^{(i)}}{\ell p_k^{(i)}} \frac{\partial p_k^{(i)}}{\partial e_3^{(i)}} n_1^{(i)} \\ &= \sum_{i \in \{1, \dots, \ell\}, z_3^{(i)} > 0} \left(-\frac{t_3^{(i)}}{\ell p_3^{(i)}} (1 - p_3^{(i)}) - \sum_{k=1}^2 \frac{t_k^{(i)}}{\ell p_k^{(i)}} (-p_k^{(i)} p_3^{(i)}) \right) n_1^{(i)} \\ &= \frac{1}{\ell} \sum_{i \in \{1, \dots, \ell\}, z_3^{(i)} > 0} \left(-t_3^{(i)} (1 - p_3^{(i)}) + \sum_{k=1}^2 t_k^{(i)} p_3^{(i)} \right) n_1^{(i)} \\ &= \frac{1}{\ell} \sum_{i \in \{1, \dots, \ell\}, z_3^{(i)} > 0} \left(-t_3^{(i)} + p_3^{(i)} \sum_{k=1}^3 t_k^{(i)} \right) n_1^{(i)} = \frac{1}{\ell} \sum_{i \in \{1, \dots, \ell\}, z_3^{(i)} > 0} (p_3^{(i)} - t_3^{(i)}) n_1^{(i)} \end{aligned}$$