

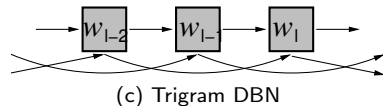
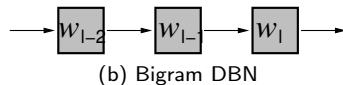
# COM4511 Speech Technology: Advanced Language Models

Anton Ragni

---

March 2, 2020





- $N$ -gram language model for general order  $n$

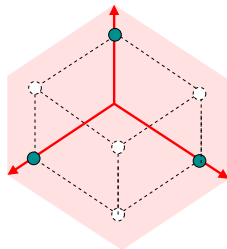
$$P(\mathbf{w}_{1:L}) = P(w_1, \dots, w_L) = \prod_{l=1}^L P(w_l | w_{l-1}, \dots, w_1) \approx \prod_{l=1}^L P(w_l | w_{l-1}, \dots, w_{l-n+1})$$

- words further than  $n - 1$  positions in the past does not affect probability
- Two primary issues
  - Markov assumption
  - discrete word representation

- ▶ Any discrete value can be mapped to continuous, vector, representation

$$\mathbf{w}_l = \begin{bmatrix} \delta(w_l, w^{(1)}) \\ \delta(w_l, w^{(2)}) \\ \vdots \\ \delta(w_l, w^{(|\mathcal{V}|)}) \end{bmatrix}$$

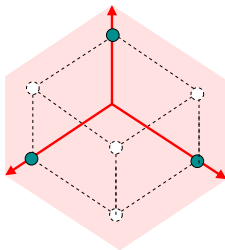
(a) Vector representation



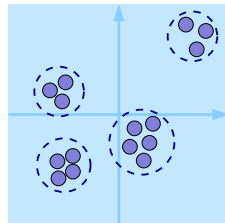
(b) Coordinate space

- ▶ each distinct discrete value is assigned a unique coordinate
- ▶ BUT does not encode any word relationship information
- ▶ Example
  - ▶ vocabulary of digits  $\mathcal{V} = \{\text{zero, one, two, three, four, five, six, seven, eight, nine}\}$
  - ▶ write down one-hot encoding matrix  $\mathbf{V}_{N \times M}$  for 21925, what are  $N$  and  $M$ ?

- Project one-hot vectors into a lower-dimensional space



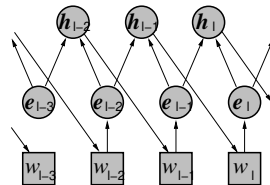
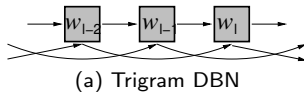
(a) 1-of- $K$  space



- Projected vector
- Semantic "equivalence"

(b) Embedding space

- optimal embedding is space where location and distances reflect word semantics
- Multiple criteria possible
  - co-occurrence statistics, word2vec, language model



(b) FFNN Trigram Pseudo DBN

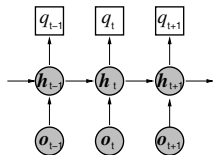
- Use feed-forward neural network to predict  $n$ -gram probabilities
  - project 1-of- $K$  vectors into a common embedding space

$$\mathbf{e}_{l-1} = \mathbf{E}_{N_1 \times |\mathcal{V}|} \mathbf{w}_{l-1_{|\mathcal{V}| \times 1}}, \quad \text{where } N_1 \ll |\mathcal{V}|$$

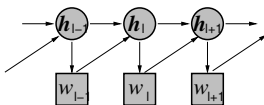
- add suitable number of feed-forward layers  $\mathbf{h}^{(l)} = \phi^{(l)}(\mathbf{A}^{(l)} \mathbf{h}^{(l-1)} + \mathbf{b}^{(l)})$
- softmax non-linearity at the final year yields

$$P(w_l | w_{l-1}, \dots, w_{l-n+1}) \triangleq h_{w_l}^{(L)}$$

- Advantages compared to maximum likelihood estimates
  - continuous word representation
  - yields probabilities for all, seen and unseen,  $n$ -grams



(a) Acoustic model  $P(\mathbf{q}_{1:T} | \mathbf{O}_{1:T})$



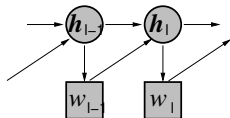
(b) Language model  $P(\mathbf{w}_{1:L})$

- Relax Markov assumption by modelling complete word history

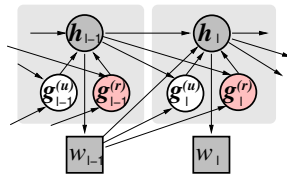
$$P(\mathbf{w}_{1:L}) = \prod_{l=1}^L P(w_l | w_{l-1}, \dots, w_1) \approx \prod_{l=1}^L P(w_l | w_{l-1}, \mathbf{h}_{l-1}) \approx \prod_{l=1}^L P(w_l | \mathbf{h}_l)$$

- a theoretically optimal language model given ... (name conditions)

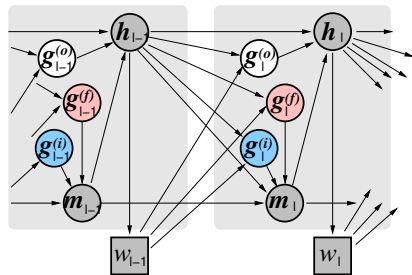
# Recurrent Units (Revisited)



(a) Simple recurrent unit



(b) Gated recurrent unit



(b) Long short-term memory

## ► Simple recurrent unit

$$\mathbf{h}_l = \phi(\mathbf{A}^{(h)} \mathbf{h}_{l-1} + \mathbf{C}^{(h)} \mathbf{w}_{l-1} + \mathbf{b}^{(h)})$$

- parameters  $\mathbf{A}_{H \times H}^{(h)}$ ,  $\mathbf{C}_{H \times |\mathcal{V}|}^{(h)}$  and  $\mathbf{b}_{H \times 1}^{(h)}$ , one-hot encoding  $\mathbf{w}_{l-1}$ , non-linearity  $\phi$  (sigmoid)
- all words embedded into continuous space parameterised by  $\mathbf{C}^{(h)}$

- ▶ Standard criterion for recurrent neural networks is **cross-entropy**\*

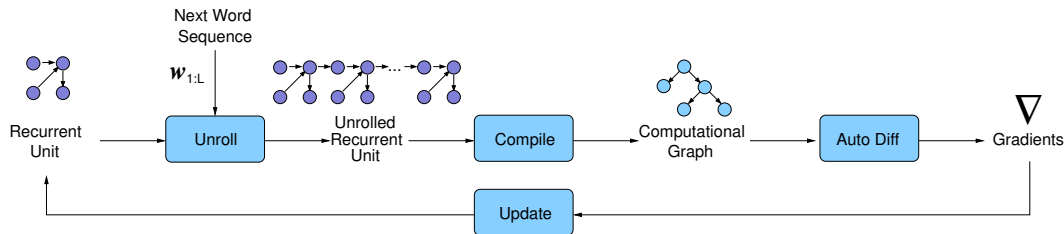
$$\mathcal{H}(P, Q) = - \sum_{\mathbf{w}} P(\mathbf{w}) \log(Q(\mathbf{w})) \approx - \frac{1}{R} \sum_{r=1}^R \log(Q(\mathbf{w}_{1:L_r}^{(r)}; \theta)) \triangleq \mathcal{L}(\mathcal{D}; \theta)$$

- ▶ use empirical estimate since true distribution  $P(\mathbf{w})$  unknown
- ▶ Use stochastic approximation to simplify parameter estimation
  - ▶ sample a minibatch of  $M$  sequences
  - ▶ update parameters
  - ▶ repeat till meet termination conditions
- ▶ Standard termination conditions
  - ▶ fixed number of iterations
  - ▶ increase in cross-entropy on held-out validation set

\* Also called negative log-likelihood criterion

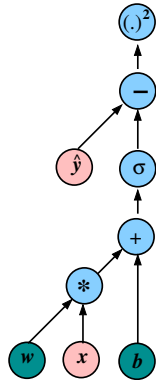


# Back-Propagation Through Time (BPTT)

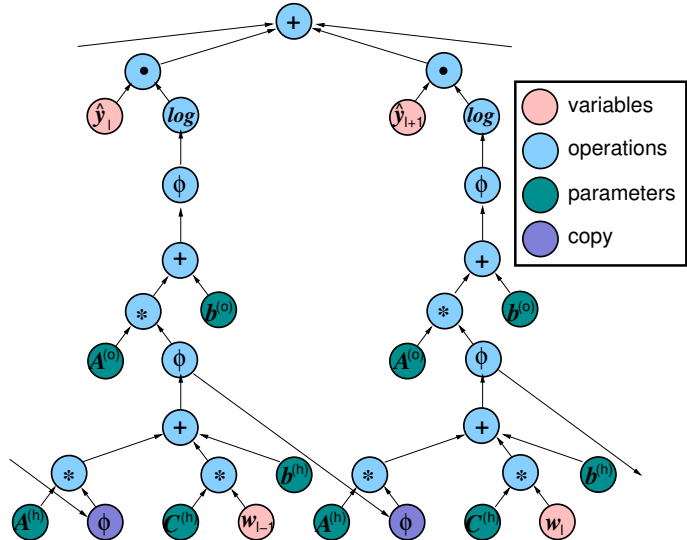


- ▶ Unroll recurrent unit to desired length
  - ▶ options available: 1-step,  $N$ -steps,  $L$ -steps where  $N < L$
  - ▶ parameters are tied across time steps
- ▶ Compute gradient with respect to recurrent unit parameters
  - ▶ compile unrolled recurrent unit into a computational graph
  - ▶ use automatic differentiation to compute gradients
- ▶ Update recurrent unit and output distribution parameters
  - ▶ and repeat!

# Automatic Differentiation: Computational Graph

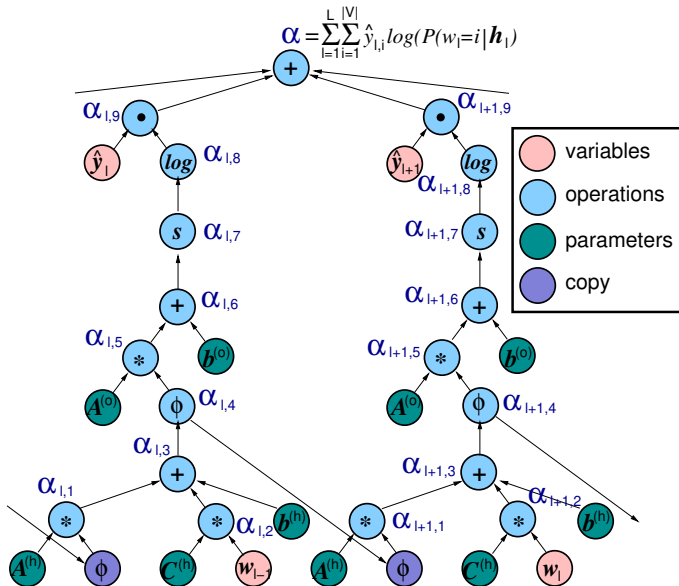


(a) Graph from past lecture



(b) Graph for simple recurrent unit

# Automatic Differentiation: Forward Pass



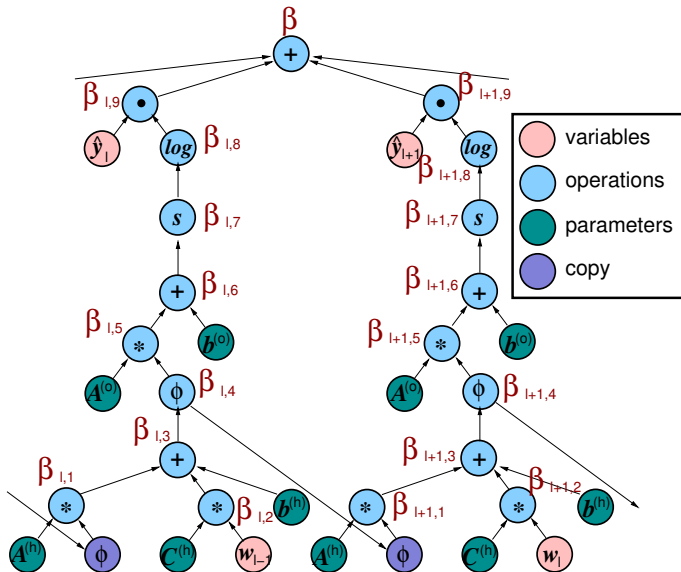
► Compute values:

$$\begin{aligned} \alpha_{l,1} &= \mathbf{A}^{(h)} * \alpha_{l-1,4} \\ \alpha_{l,2} &= \mathbf{C}^{(h)} * \mathbf{w}_{l-1} \\ \alpha_{l,3} &= \alpha_{l,1} + \alpha_{l,2} + \mathbf{b}^{(h)} \\ \alpha_{l,4} &= \phi(\alpha_{l,3}) \\ \alpha_{l,5} &= \mathbf{A}^{(o)} * \alpha_{l,4} \\ \alpha_{l,6} &= \alpha_{l,5} + \mathbf{b}^{(o)} \\ \alpha_{l,7} &= s(\alpha_{l,6}) \\ \alpha_{l,8} &= \log(\alpha_{l,7}) \\ \alpha_{l,9} &= \hat{y}_l \cdot \alpha_{l,8} \\ \alpha &= \alpha + \alpha_{l,9} \end{aligned}$$

where  $s$  is softmax,  $\phi$  is sigmoid

$$\hat{y}_{l,i} = \begin{cases} 1, & \text{if } w_l^{(i)} = w_l \\ 0, & \text{otherwise} \end{cases}$$

# Automatic Differentiation: Backward Pass



► Compute derivatives:

$$\begin{aligned}\beta &= \frac{\partial \alpha}{\partial \alpha} \equiv 1 \\ \beta_{l,9} &= \beta \frac{\partial \alpha}{\partial \alpha_{l,9}} \\ \beta_{l,8} &= \beta_{l,9} \frac{\partial \alpha_{l,9}}{\partial \alpha_{l,8}} \\ \beta_{l,7} &= \beta_{l,8} \frac{\partial \alpha_{l,8}}{\partial \alpha_{l,7}} \\ \beta_{l,6} &= \beta_{l,7} \frac{\partial \alpha_{l,7}}{\partial \alpha_{l,6}} \\ \beta_{l,5} &= \beta_{l,6} \frac{\partial \alpha_{l,6}}{\partial \alpha_{l,5}} \\ \beta_{l,4} &= \beta_{l,5} \frac{\partial \alpha_{l,5}}{\partial \alpha_{l,4}} \\ \beta_{l,3} &= \beta_{l,4} \frac{\partial \alpha_{l,4}}{\partial \alpha_{l,3}} \\ \beta_{l,2} &= \beta_{l,3} \frac{\partial \alpha_{l,3}}{\partial \alpha_{l,2}} \\ \beta_{l,1} &= \beta_{l,3} \frac{\partial \alpha_{l,3}}{\partial \alpha_{l,1}}\end{aligned}$$

► Accumulate derivatives:

$$\frac{\partial \alpha}{\partial \mathbf{A}^{(h)}} = \frac{\partial \alpha}{\partial \mathbf{A}^{(h)}} + \beta_{l,1} \frac{\partial \alpha_{l,1}}{\partial \mathbf{A}^{(h)}}$$

$$\frac{\partial \alpha}{\partial \mathbf{C}^{(h)}} = \frac{\partial \alpha}{\partial \mathbf{C}^{(h)}} + \beta_{l,2} \frac{\partial \alpha_{l,2}}{\partial \mathbf{C}^{(h)}}$$

$$\frac{\partial \alpha}{\partial \mathbf{b}^{(h)}} = \frac{\partial \alpha}{\partial \mathbf{b}^{(h)}} + \beta_{l,3} \frac{\partial \alpha_{l,3}}{\partial \mathbf{b}^{(h)}}$$

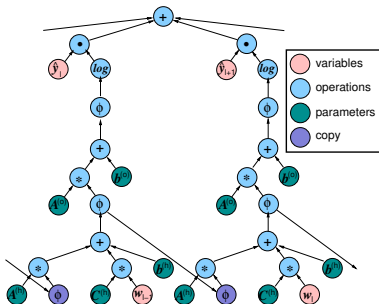
$$\frac{\partial \alpha}{\partial \mathbf{A}^{(o)}} = \frac{\partial \alpha}{\partial \mathbf{A}^{(o)}} + \beta_{l,5} \frac{\partial \alpha_{l,5}}{\partial \mathbf{A}^{(o)}}$$

$$\frac{\partial \alpha}{\partial \mathbf{b}^{(o)}} = \frac{\partial \alpha}{\partial \mathbf{b}^{(o)}} + \beta_{l,6} \frac{\partial \alpha_{l,6}}{\partial \mathbf{b}^{(o)}}$$

► Update parameters:

$$\hat{\mathbf{A}}^{(h)} = \mathbf{A}^{(h)} + \rho \frac{\partial \alpha}{\partial \mathbf{A}^{(h)}}$$

- similar rules for other parameters
- why derivative is added?



- ▶ Algebraic manipulations with high-dimensional data expensive
  - ▶ graphics processing units (GPU) provide support for efficient computation
  - ▶ BUT operations needs to be arranged efficiently
- ▶ Options available to construct mini-batches

<s>	all	you	get	is	risk	<s>	brokers	...	don't
deny	that	<s>	heating	oil	prices	rose	<s>	...	but
will	it	be	enough	<s>	the	lobby	says	...	<s>

(a) Cropping

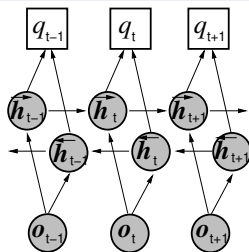
<s>	how	interesting	<s>	€	€	€
<s>	all	you	get	is	risk	<s>
<s>	don't	deny	that	<s>	€	€

(b) Sentence bunching

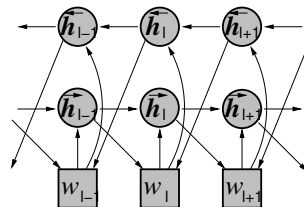
<s>	how	interesting	<s>	all	you	...	need	that	<s>	€
<s>	the	lobby	says	no	<s>	...	enough	<s>	€	€

(c) Sentence splicing

- ▶ only last two options yield a valid language model (why?)



(a) Acoustic model

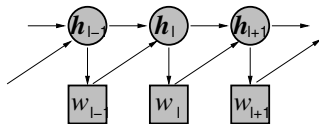


(a) Language model

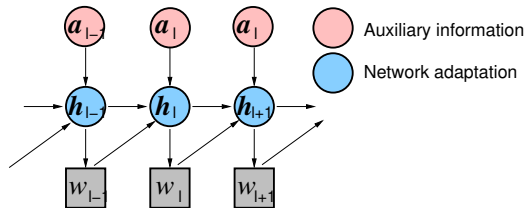
- Use information from **future words** to help predicting current words

$$P(\mathbf{w}_{1:L}) \triangleq \frac{1}{Z} \tilde{P}(\mathbf{w}_{1:L}) = \frac{1}{Z} \prod_{l=1}^L P(w_l | \mathbf{w}_{1:l-1}, \mathbf{w}_{l+1:L}) \approx \frac{1}{Z} \prod_{l=1}^L P(w_l | \vec{\mathbf{h}}_{l-1}, \overleftarrow{\mathbf{h}}_{l+1})$$

- compensate for inadequate assumptions/approximations of unidirectional models
  - normalisation term  $Z = \sum_{\mathbf{w}} \tilde{P}(\mathbf{w})$  cannot be efficiently computed (exc trivial cases)
- Performance assessed using **pseudo-perplexity** based on unnormalised probabilities
  - **discuss possible issues**
- **Exercise:** draw pseudo-DBNs for GRU and LSTM



(a) Standard language model



(b) Adaptation options

- ▶ Incorporate auxiliary information
  - ▶ topic and genre adaptation
  - ▶ BUT need a suitable representation
- ▶ Adjust parameters and/or activation functions
  - ▶ fine-tune model and/or activation parameters on adaptation data
  - ▶ need adaptation data
- ▶ Other options possible



- ▶ This lecture examined advanced language models
  - ▶ continuous word representations
  - ▶ full word history modelling
- ▶ Focus on recurrent forms of language models
  - ▶ recurrent units
  - ▶ parameter estimation
  - ▶ future words
  - ▶ adaptation
- ▶ Not examined in this lecture
  - ▶ convolutional language models