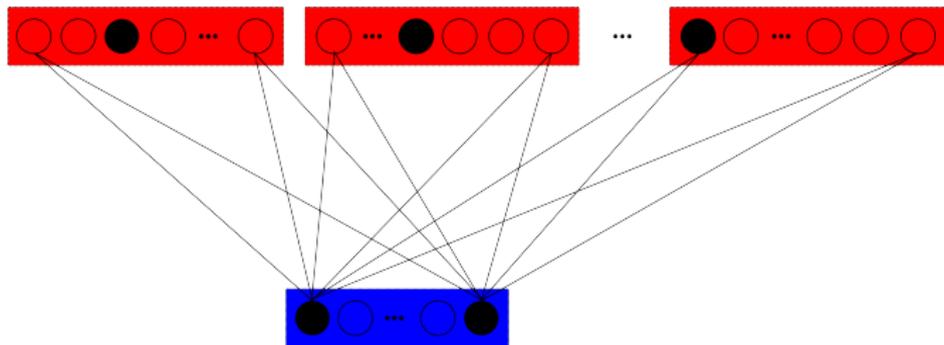


# Continuous Space Language Models using Restricted Boltzmann Machines

Jan Niehues and Alex Waibel

Institute for Anthropomatics



# Motivation

## N-gram based language models

- Make use of large corpora
- Can be trained efficiently

## Domain Adaptation

- Language models trained on small corpora are needed
- Language model has to back-off to smaller contexts
- Continuous space language models always use same context size
- Longer training time not as problematic
- Aim: Application during decoding

## Related work

First approaches of neural network language models using word categories in the 90s (Nakamura et al., 1990)

CSLM for speech recognition (Bengio et al., 2003, Schwenk et al., 2002)

SOUL Language model by Le et al., 2011

RBM-based language model (Mnih and Hinton, 2007)

# Overview

Restricted Boltzmann Machine

RBM for Language Model

- Calculating probabilities efficiently

Evaluation

- German-English
- English-French

Conclusion

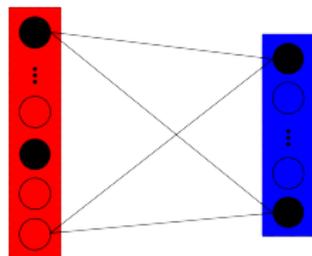
# Restricted Boltzmann Machine

## Layout

- 2 layer neuronal net
- Binary units
- Weighted connections between the layers
- No connections within the layer

input layer

hidden layer



# Restricted Boltzmann Machine

## Layout

## Probability

- Probability defined by the energy

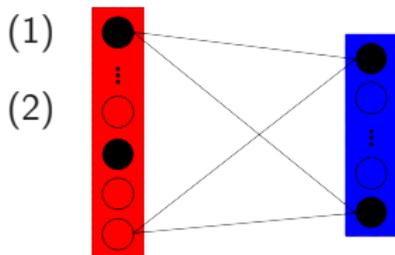
$$p(v, h) = \frac{1}{Z} e^{-E(v, h)}$$

$$E(v, h) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i, j} v_i h_j w_{ij}$$

- Problem: Hidden state needs to be known

input layer

hidden layer



## Layout

## Probability

- Probability defined by the energy

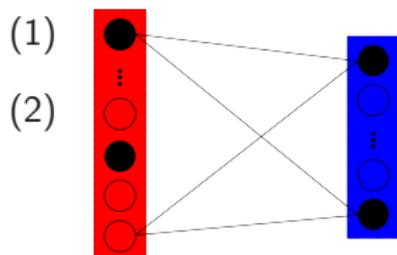
$$p(v, h) = \frac{1}{Z} e^{-E(v, h)}$$

$$E(v, h) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i, j} v_i h_j w_{ij}$$

- Problem: Hidden state needs to be known

input layer

hidden layer



## Layout

## Probability

- Probability defined by the energy
- Probability using free energy

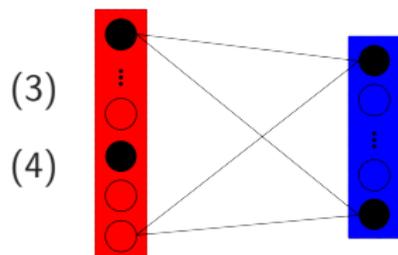
$$p(v) = \frac{1}{Z} e^{-F(v)}$$

$$F(v) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{j \in \text{hidden}} \log(1 + e^{x_j})$$

$$x_j = b_j + \sum_{i \in \text{visible}} v_i w_{ij} \quad (5)$$

input layer

hidden layer



# Restricted Boltzmann Machine

Layout

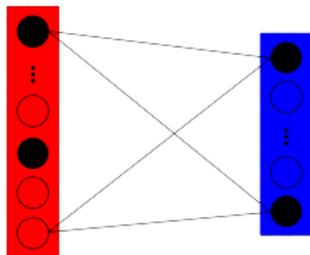
Probability

Training

- Contrastive Divergence
- Increase probability of seen training example

input layer

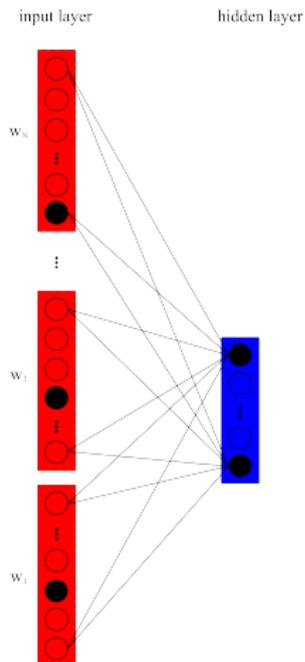
hidden layer



# RBM for Language modeling

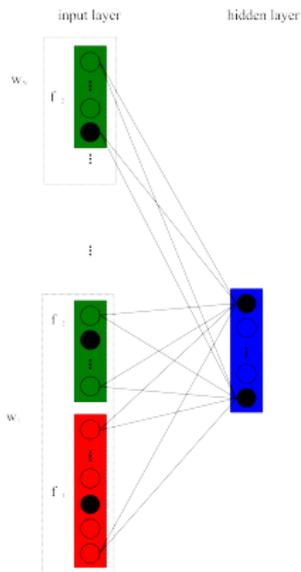
## Layout

- N input blocks
- Each of the block has V units
- H hidden units
- $N*V*H$  weights and  $N*V + H$  biases



## Layout

- N input blocks
- Each of the block has V units
- H hidden units
- $N \cdot V \cdot H$  weights and  $N \cdot V + H$  biases
- Easy integration of additional word factors
- Replace block by W sub blocks

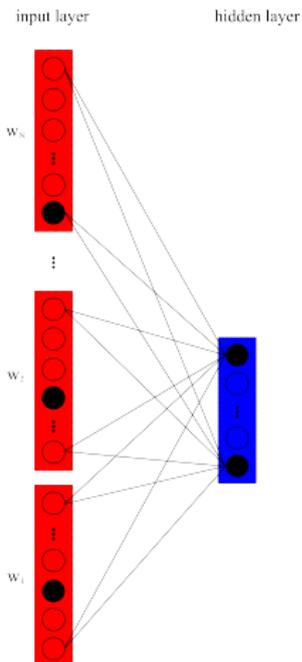


# RBM for Language modeling

## Layout

## Training

- Contrastive Divergence
- Random order of the n-grams in training data
- 1 iteration over all the data
- 1 iteration of Gibbs sampling for collection samples



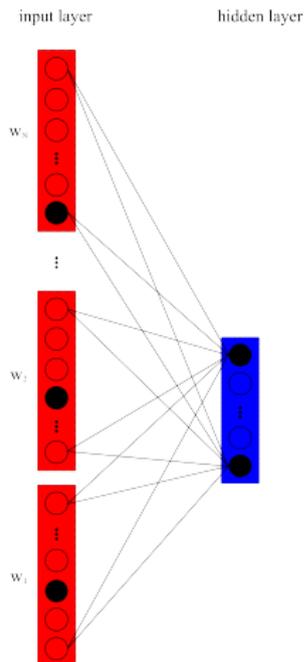
# N-gram Probability

Use language model in the decoder

→ efficient calculation needed

- Use free energy instead of probability
  - No normalization needed
- Complexity

$$\begin{aligned}
 F(v) &= - \sum_{i \in \text{visible}} a_i v_i \\
 &\quad - \sum_{j \in \text{hidden}} \log(1 + e^{x_j}) \\
 x_j &= b_j + \sum_{i \in \text{visible}} v_i w_{ij}
 \end{aligned}$$



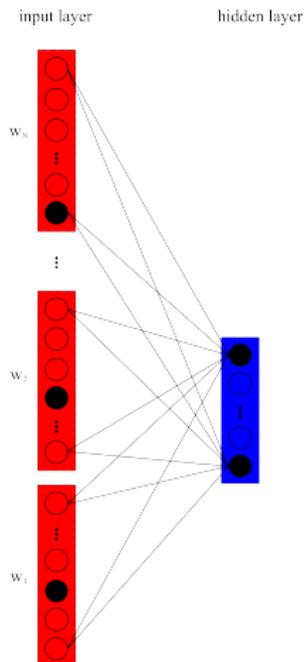
# N-gram Probability

Use language model in the decoder  
 → efficient calculation needed

- Use free energy instead of probability
  - No normalization needed
- Complexity

$$F(v) = -O(N) - \sum_{j \in \text{hidden}} \log(1 + e^{x_j})$$

$$x_j = b_j + O(N)$$

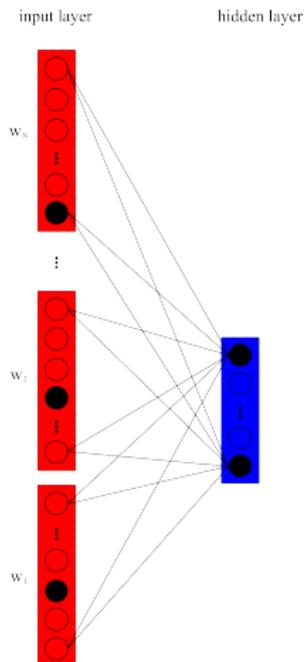


# N-gram Probability

Use language model in the decoder  
→ efficient calculation needed

- Use free energy instead of probability
  - No normalization needed
- Complexity

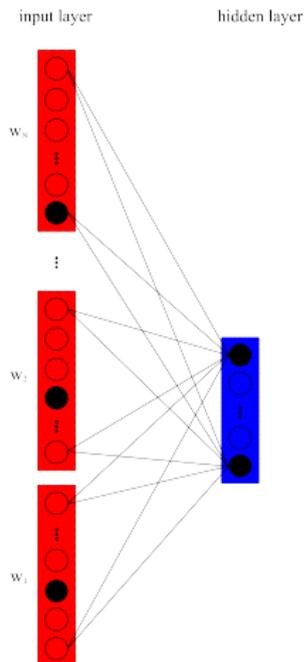
$$F(v) = -O(N) \\ - O(H * N) \\ x_j = b_j + O(N)$$



# N-gram Probability

Use language model in the decoder  
 → efficient calculation needed

- Use free energy instead of probability
  - No normalization needed
- Calculate free energy in  $O(H * N)$
- Independent of the vocabulary size



Feature describing probability of the whole sentence

$$\sum_{j \in L+N-1} F(w_{j-N+1} \dots w_j) \quad (6)$$

- Sum over free energy of all n-grams

<s>	<s>	<s>			go	home	</s>
<s>	<s>		go	<b>go</b>	<b>home</b>	</s>	</s>
<s>		go	home	<b>home</b>	</s>	</s>	</s>

Feature describing probability of the whole sentence

$$\sum_{j \in L+N-1} F(w_{j-N+1} \dots w_j) \quad (7)$$

- Sum over free energy of all n-grams
- proportional to approximation of the geometric mean of all language model probabilities with contexts  $\leq N$ 
  - terms depending on the sentence length are ignored
- easy to integrate into the decoder

## System description

### German to English

- TED Translation task
- Phrase-based SMT system
- Training data: EPPS, NC, TED, BTEC
- In-domain: TED
- Additional word factor: automatic word classes generated by MKCLS ( 50 classes)

### English to French

- System built during IWSLT 2012 Evaluation

## German to English

Table: *Experiments on German to English*

System	Iterations	BLEU Score	
		Dev	Test
Baseline		26.31	23.02
+ RBMLM H32	1	27.39	23.82
+ RBMLM H32	10	27.61	24.47
+ FRBMLM H32	1	27.54	24.15
Baseline+NGRAM		27.45	24.06
+ RBMLM H32	1	27.64	24.33
+ RBMLM H32	10	27.95	24.38
+ FRBMLM H32	1	27.80	24.40

## Number of hidden units

Table: *Experiments using different number of hidden units*

System	Hidden Units	BLEU Score	
		Dev	Test
NGRAM		27.09	23.80
RBMLM	8	25.65	23.16
	16	25.67	23.07
	32	26.40	23.41
	64	26.12	23.18

## Training iterations

Table: *Experiments using different number of training iterations*

System	Iterations	No Large LM		Large LM	
		Dev	Test	Dev	Test
NGRAM		27.09	23.80	27.45	24.06
RBMLM	1	26.40	23.41	27.39	23.82
	5	26.72	23.38	27.40	23.98
	10	26.90	23.51	27.61	24.47
	15	26.57	23.47	27.63	24.22
	20	26.16	23.20	27.49	24.30

## English to French

Table: *Experiments on English to French*

System	BLEU Score	
	Dev	Test
Baseline	28.93	31.90
RBMLM	28.99	31.76
FRBMLM	29.02	32.03

- Continuous space language model using Restricted Boltzmann Machines
- Approximations to efficiently calculate language model score
- Language model score is independent of vocabulary size
- Integration into decoding
- Factor language model
- Experiments on two TED translation tasks
  - Detailed experiments on German-English
  - Slight improvements on English-French

## Example

$$P_3(S) = P(I | \langle s \rangle \langle s \rangle) * P(\text{go} | \langle s \rangle I) * P(\text{home} | I \text{go}) * P(\langle /s \rangle | \text{gohome})$$

$$P_2(S) = P(I | \langle s \rangle) * P(\text{go} | I) * P(\text{home} | \text{go}) * P(\langle /s \rangle | \text{home})$$

$$P_1(S) = P(I) * P(\text{go}) * P(\text{home})$$

$$\begin{aligned} P(S) &= \sqrt[3]{P_3(S) * P_2(S) * P_1(S)} \\ &= \frac{P(\langle s \rangle \langle s \rangle I)}{P(\langle s \rangle \langle s \rangle)} * \frac{P(\langle s \rangle I)}{P(\langle s \rangle)} * P(I) \\ &\quad * \frac{P(\langle s \rangle I \text{go})}{P(\langle s \rangle I)} * \frac{P(I \text{go})}{P(I)} * P(\text{go}) \\ &\quad * \frac{P(I \text{gohome})}{P(I \text{go})} * \frac{P(\text{gohome})}{P(\text{go})} * P(\text{home}) \\ &\quad * \frac{P(\text{gohome} \langle /s \rangle)}{P(\text{gohome})} * \frac{P(\text{home} \langle /s \rangle)}{P(\text{home})} \end{aligned}$$

## Example

$$\begin{aligned} &= \frac{P(\langle s \rangle \langle s \rangle I)}{P(\langle s \rangle \langle s \rangle)} * \frac{P(\langle s \rangle I)}{P(\langle s \rangle)} * P(I) \\ &* \frac{P(\langle s \rangle Igo)}{P(\langle s \rangle I)} * \frac{P(Igo)}{P(I)} * P(go) \\ &* \frac{P(Igohome)}{P(Igo)} * \frac{P(gohome)}{P(go)} * P(home) \\ &* \frac{P(gohome \langle /s \rangle)}{P(gohome)} * \frac{P(home \langle /s \rangle)}{P(home)} \\ &= P(\langle s \rangle \langle s \rangle I) P(\langle s \rangle Igo) * P(Igohome) \\ &* P(gohome \langle /s \rangle) * P(home \langle /s \rangle) \\ &= P(\langle s \rangle \langle s \rangle I) P(\langle s \rangle Igo) * P(Igohome) \\ &* P(gohome \langle /s \rangle) P(home \langle /s \rangle \langle /s \rangle) \end{aligned}$$