

# DEEP LEARNING FOR SPEECH & LANGUAGE

Winter Seminar UPC TelecomBCN, 24 - 31 January 2017



## Instructors



Antonio  
Bonafonte

J. Adrián Rodríguez  
Fonollosa

Marta R.  
Costa-jussà

Javier  
Hernando

Santiago  
Pascual

Elisa  
Sayrol

Xavier  
Giró

## Organizers



Image Processing Group  
Signal Theory and Communications Department



UNIVERSITAT POLITÈCNICA  
DE CATALUNYA  
BARCELONATECH

+ info: [TelecomBCN.DeepLearning.Barcelona](https://www.telecombcn.com/deeplearning-barcelona)

[\[course site\]](#)

Day 3 Lecture 1

# Language Model

Marta R. Costa-jussà

# Previous concepts from this course

- Word embeddings
- Feed-forward network and softmax
- Recurrent neural network (handle variable-length sequences)

# What is the most probable sentence?

Two birds are flying

Two beards are flying

## Probability of a sentence

- Suppose you record a database of one billion utterances in English.
- If the sentence “how's it going?” appears 76,413 times in that database, then we say  
— $P(\text{how's it going?}) = 76,413 / 1,000,000,000$

# A language model finds the probability of a sentence

- Given a sentence  $(w_1, w_2, \dots, w_T)$ ,
- What is  $p(w_1, w_2, \dots, w_T) = ?$

# **An n-gram language model**

# Chain rule probability and Markov simplifying assumption

$$\frac{p(w_1, w_2, \dots, w_T)}{p(w_1)} = p(w_T | w(T-1), w(T-2) \dots w_1) p(w(T-1) | w(T-2), w(T-3) \dots w_1) \dots$$

**Markov simplifying assumption:** The current word only depends on  $n$  previous words.

$$p(w_t | w(t-1)w(t-2) \dots w_1) \sim p(w_t | w(t-1))$$

## Objective

$$\begin{aligned} p(w_1, w_2, \dots w_T) &= \prod_{t=1}^T p(w_t | w_1 \dots w_{t-1}) \\ &\approx \prod_{t=1}^T p(w_t | w_{t-1} \dots w_{t-n}) \end{aligned}$$



# An n-gram-based language model

- An n-word substring is called an n-gram.
- If  $n=1$ , we say unigram; if  $n=2$ , we say bigram; if  $n=3$ , we say trigram.
- $P(< s> \text{ I like snakes that are not poisonous } </s>) \sim$   
 $b(\text{I} \mid < s>) * b(\text{like} \mid \text{I}) * b(\text{snakes} \mid \text{like}) * \dots * b(\text{poisonous} \mid \text{not})$   
 $* b(</s> \mid \text{poisonous})$

# An n-gram-based language model

Unigram probabilities

$$p(w_1) = \frac{\text{count}(w_1)}{\text{total words observed}}$$

Bigram probabilities

$$p(w_2|w_1) = \frac{\text{count}(w_1w_2)}{\text{count}(w_1)}$$

Trigram probabilities

$$p(w_3|w_1w_2) = \frac{\text{count}(w_1w_2w_3)}{\text{count}(w_1w_2)}$$

**Any issues with the above model?**

# Some examples...

"<s> que la fuerza te acompañe </s>", = *may the force be with you*

bigrams and trigrams like:

*fuerza te*

*la fuerza te*

*fuerza te acompañe*

*te acompañe </s>*

do not appear in the big corpus of El Periodico (40 M words)

BUT PROBABILITY OF THE SENTENCE SHOULD NOT BE ZERO!!!!

# Sparse counts are a big problem

- Backing off avoids zero probabilities

$$.8 * p(w_3 | w_1 w_2)$$

$$+.15 * p(w_3 | w_2)$$

$$+0.049 * p(w_3)$$

$$+.001$$

# Sparse counts are a big problem

- Smoothing avoids zero probabilities

$$.8 * p(w_3 | w_1 w_2)$$

$$+.15 * p(w_3 | w_2)$$

$$+0.049 * p(w_3)$$

$$+.001$$

**Any other issue?**

# Lack of generalization

Mary buys two apples and two oranges in the market  
three apples are for me

the tree has three oranges

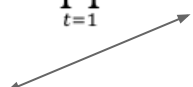


# A neural language model

To generalize to un-seen n-grams

# A neural language model

Find a function that takes as input  $n-1$  words and returns a conditional probability of the next one

$$p(w_1, w_2, \dots, w_T) = \prod_{t=1}^T p(w_t | w_1 \dots w_{t-1})$$
$$\approx \prod_{t=1}^T p(w_t | w_{t-1} \dots w_{t-n})$$

$$p(w_t | w_{t-n}, \dots, w_{t-1}) = f_{w_t}(w_{t-n}, \dots, w_{t-1})$$

# Architecture: neural language model

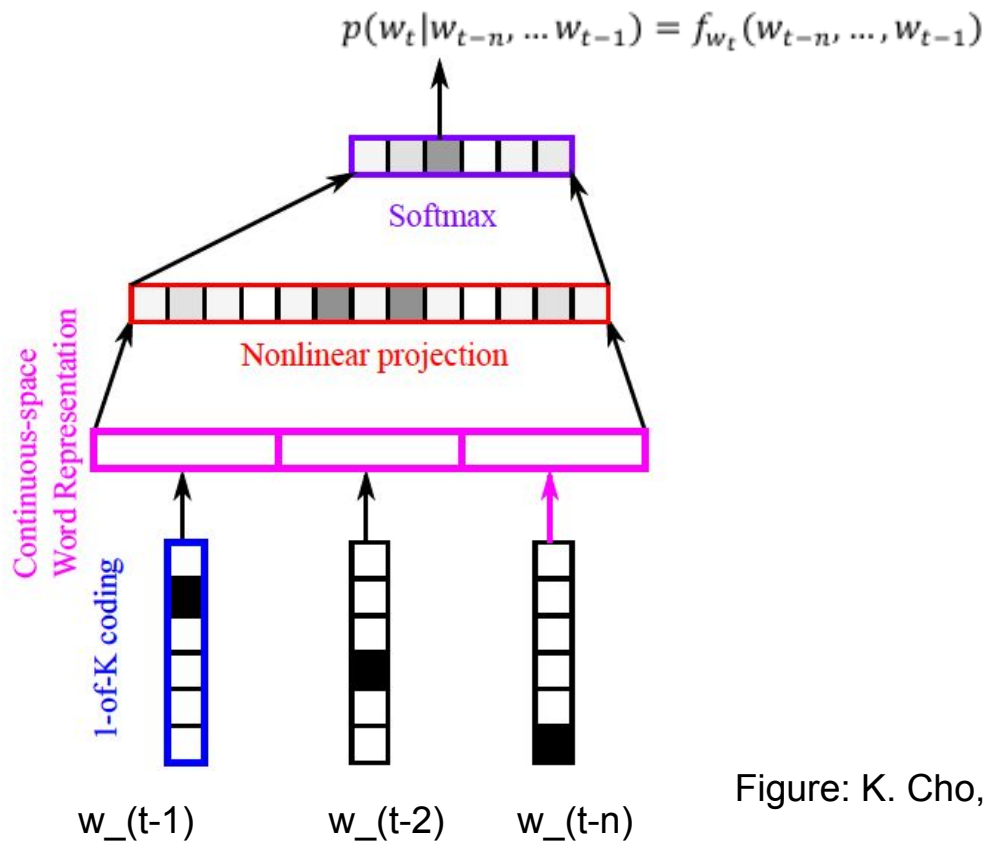
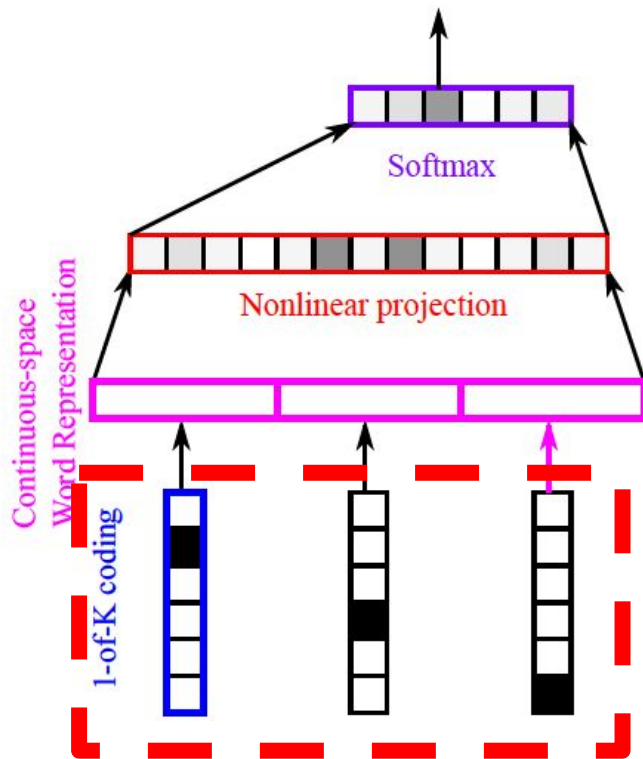


Figure: K. Cho, DL4MT course, 2015

# Architecture: representation of input words



our goal is to put the least amount of prior knowledge

Figure: K. Cho, DL4MT course, 2015

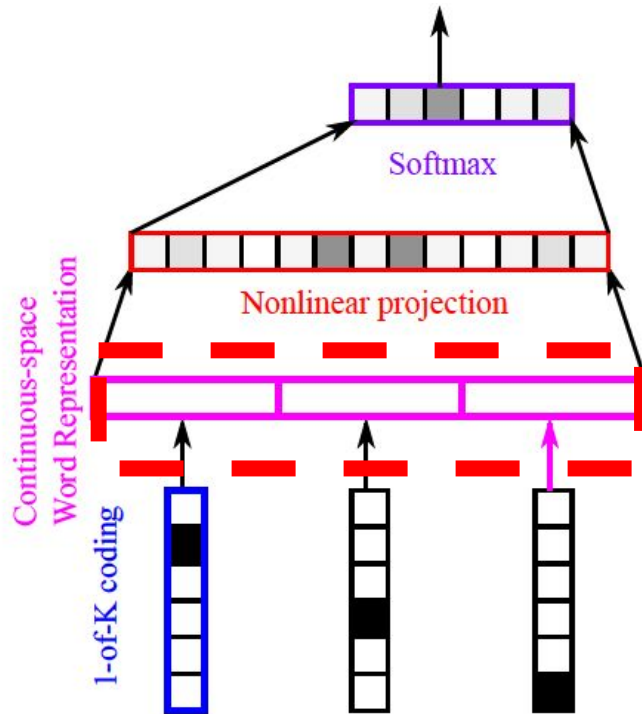
# Step 1: One-hot encoding

From previous lectures

Natural language words can be one-hot encoded on a vector of dimensionality equal to the size of the dictionary ( $K=|V|$ ).

Word	One-hot encoding
economic	000010...
growth	001000...
has	100000...
slowed	000001...

# Architecture: continuous word representation



input vectors are multiplied by the weight matrix ( $E$ ), to obtain continuous vectors

this weight matrix ( $E$ ) is also called word embedding and should reflect the meaning of a word

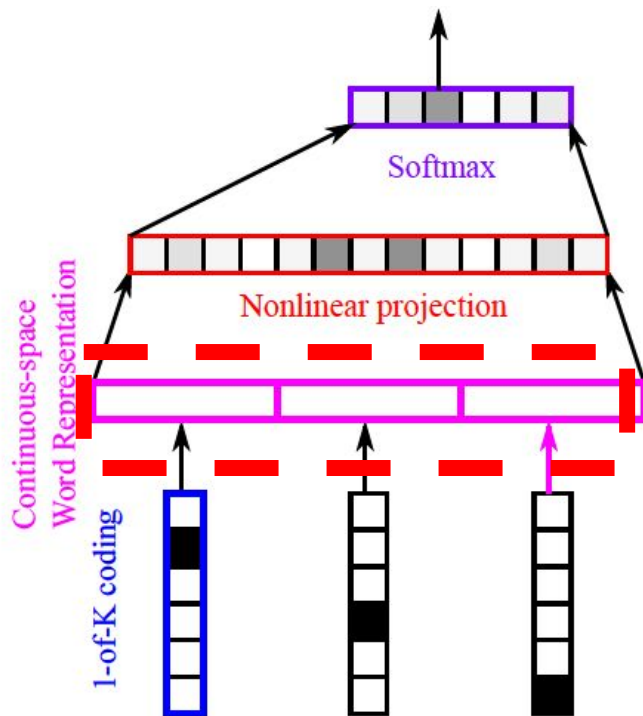
$$E = \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_{|V|} \end{bmatrix}, \quad \mathbf{e}_i \in \mathbb{R}^d.$$

$$E \in \mathbb{R}^{|V| \times d}.$$

$$E^T \mathbf{w}_i = \mathbf{e}_i.$$

Figure: K. Cho, DL4MT course, 2015

# Architecture: context vector



we get a sequence of continuous vectors, by concatenating the continuous representations of the input words

$$\mathbf{p}^j = \mathbf{E}^\top \mathbf{w}^j$$

$$\mathbf{p} = [\mathbf{p}^1; \mathbf{p}^2; \dots; \mathbf{p}^{n-1}]^\top$$

↓  
context vector

Figure: K. Cho, DL4MT course, 2015

# Architecture: nonlinear projection

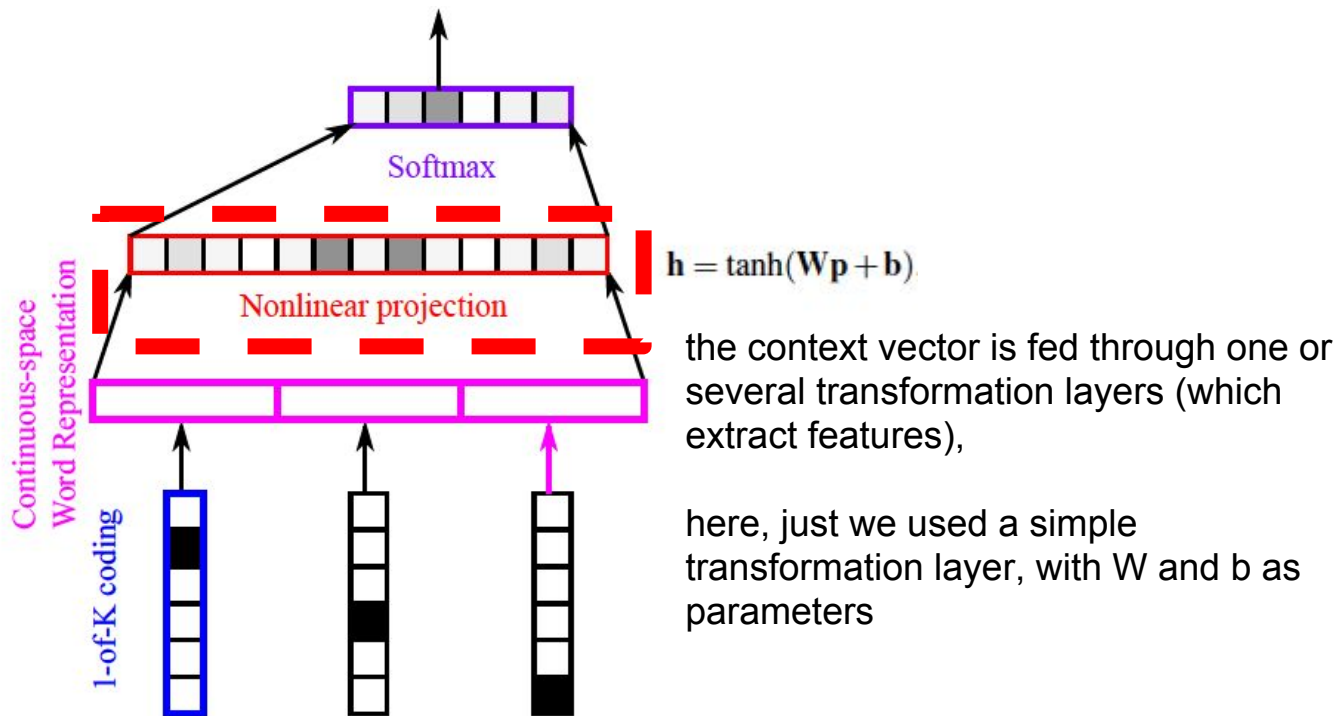


Figure: K. Cho, DL4MT course, 2015



# Architecture: output probability distribution

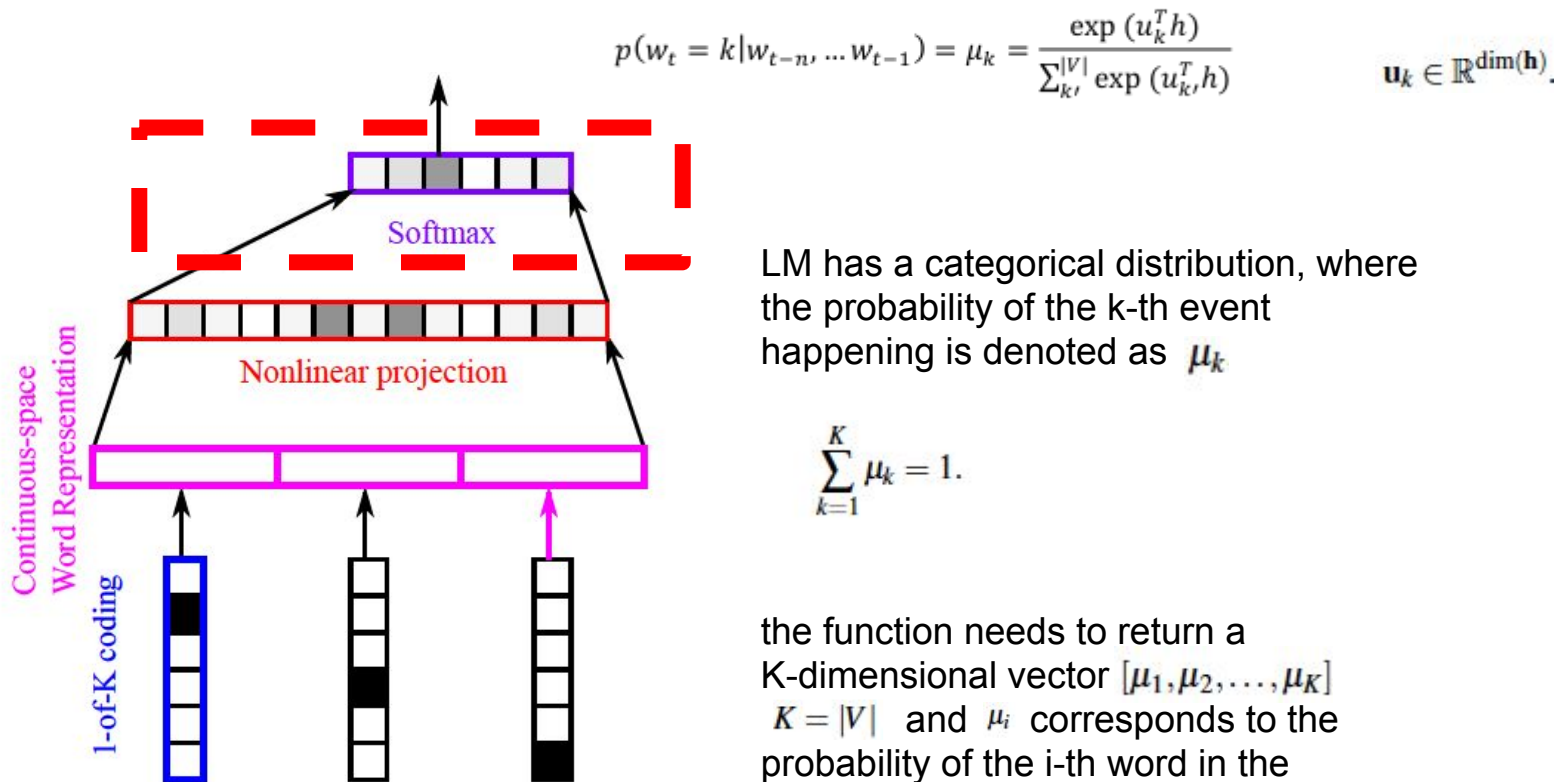


Figure: K. Cho, DL4MT course, 2015

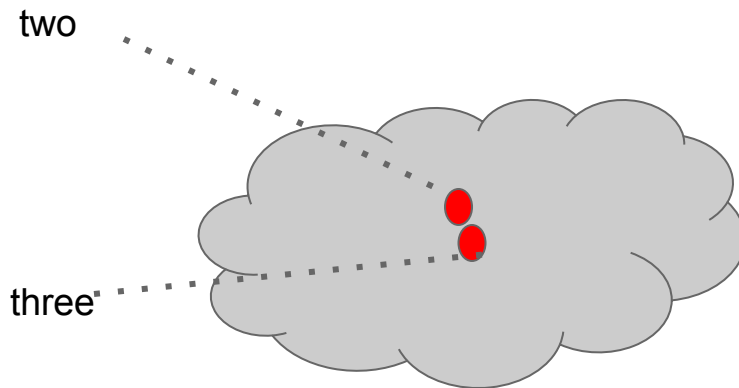
**Why this model is generalizing to  
unseen events?**

# Further generalization comes from embeddings

Mary buys two apples and two oranges in the market

three apples are for me

the tree has three oranges



# Recurrent Language Model

To further generalize to un-seen n-grams

# A neural language model

Still assumes the  $n$ -th order Markov property  
it looks only as  $n-1$  past words

In France , there are around 66 million people and they speak French.

# How we can modelate variable-length input?

$$p(w_1, w_2, \dots w_T) = \prod_{t=1}^T p(w_t | w_1 \dots w_{t-1})$$

# How we can modelate variable-length input?

$$p(w_1, w_2, \dots w_T) = \prod_{t=1}^T p(w_t | w_1 \dots w_{t-1})$$

We directly model the original conditional probabilities

Via Recursion:

$$h_0 = 0$$

Initial condition

$$h_t = f(w_{t-1}, h_{t-1})$$

Recursion



summarizes the history from  $w_1$  to  $w(t-1)$

# How we can modelate variable-length input?

$$p(w_1, w_2, \dots w_T) = \prod_{t=1}^T p(w_t | w_1 \dots w_{t-1})$$

We directly model the original conditional probabilities

Via Recursion:

$$h_0 = 0$$

Initial condition

$$h_t = f(w_{t-1}, h_{t-1})$$

Recursion



summarizes the history from  $w_1$  to  $w(t-1)$

The RNN is capable of summarizing a variable-length input sequence ( $w$ ) into a memory state ( $h$ )



# Example

$p(\text{Mary, buys, the, apple})$


- (1) Initialization:  $h_0=0 \rightarrow p(\text{Mary})=g(h_0)$
- (2) Recursion
  - (a)  $h_1=f(h_0, \text{Mary}) \rightarrow p(\text{buys}|\text{Mary})=g(h_1)$
  - (b)  $h_2=f(h_1, \text{buys}) \rightarrow p(\text{the}|\text{Mary, buys})=g(h_2)$
  - (c)  $h_3=f(h_2, \text{the}) \rightarrow p(\text{apple}|\text{Mary, buys, the})=g(h_3)$
- (3) Output:  $p(\text{Mary, buys, the, apple})=g(h_0)g(h_1)g(h_2)g(h_3)$

It works for any number of context words

READ, UPDATE, PREDICT

# A recurrent neural language model

Can modelate the probability of a variable-length input sequence

$$p(w_1, w_2, \dots w_T) = \prod_{t=1}^T p(w_t | w_1 \dots w_{t-1})$$


conditional probability that we want  
to compute

# A recurrent neural language model

what we need

(1) Transition function

$$h_t = f(w_{t-1}, h_{t-1})$$

(2) Output function

$$p(w_t = k | w_1 \dots w_{t-1})$$

# (Naive) Transition function

Inputs:

one-hot vector  $w_{(t-1)} \in \{0,1\}^{|V|}$   
+ hidden state  $h_{(t-1)} \in \mathbb{R}^d$

Parameters:

input weight matrix  $E \in \mathbb{R}^{d \times |V|}$   
+ transition weight matrix  $U \in \mathbb{R}^{d \times d}$   
+ bias vector (b)

# (Naive) Transition function

nonlinear transformation

Linear Transformation of the Previous Hidden State

$$h_t = \tanh(E^T w_{t-1} + U h_{t-1} + b)$$

Continuous-space Representation of word

The diagram shows the equation  $h_t = \tanh(E^T w_{t-1} + U h_{t-1} + b)$ . Three annotations with leader lines point to specific parts of the equation: 'nonlinear transformation' points to the  $\tanh$  function; 'Linear Transformation of the Previous Hidden State' points to the term  $U h_{t-1}$ , which is enclosed in a gray oval; and 'Continuous-space Representation of word' points to the term  $E^T w_{t-1}$ , which is also enclosed in a gray oval.

# Output function

Input:

hidden state (ht)

Parameters:

output matrix

$$V \in \mathbb{R}^{|V| \times d}$$

bias vector (c)

## Output function

This summary vector is affine-transformed followed by a softmax nonlinear function to compute the conditional probability of each word.

$$\mu = \text{softmax}(Vh_t + c)$$

**Measure to evaluate**

Perplexity



# Perplexity: a measure to evaluate language modeling

Perplexity measures how high a probability the language model assigns to correct next words in the test corpus “on average”. A better language model is the one with a lower perplexity.

Perplexity measures as well how complex is a task equals size of vocabulary (V)

$$PP=V$$

# Comparing language models

LM	Hidden Layers	PPL
n-gram-based		131.2
+feed-forward	600	112.5
+RNN	600	108.1
+LSTM	600	92.0

Results from Sundermeyer et al, 2015

# Applications of language modeling

Speech recognition

Machine translation

Handwriting recognition

Information retrieval

...

# A bad language model

HERMAN



# Summary

- Language modeling consists in assigning a probability to a sequence of words.
- We can model a sequence of words with n-grams, feed-forward networks and recurrent networks.
- Feed-forward networks are able to generalise unseen contexts
- RNN are able to use variable contexts

# Learn more

Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. 2003. A neural probabilistic language model. *J. Mach. Learn. Res.* 3 (March 2003), 1137-1155.

Martin Sundermeyer, Hermann Ney, and Ralf Schlüter. 2015. From feedforward to recurrent LSTM neural networks for language modeling. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.* 23, 3 (March 2015), 517-529.

DOI=<http://dx.doi.org/10.1109/TASLP.2015.2400218>

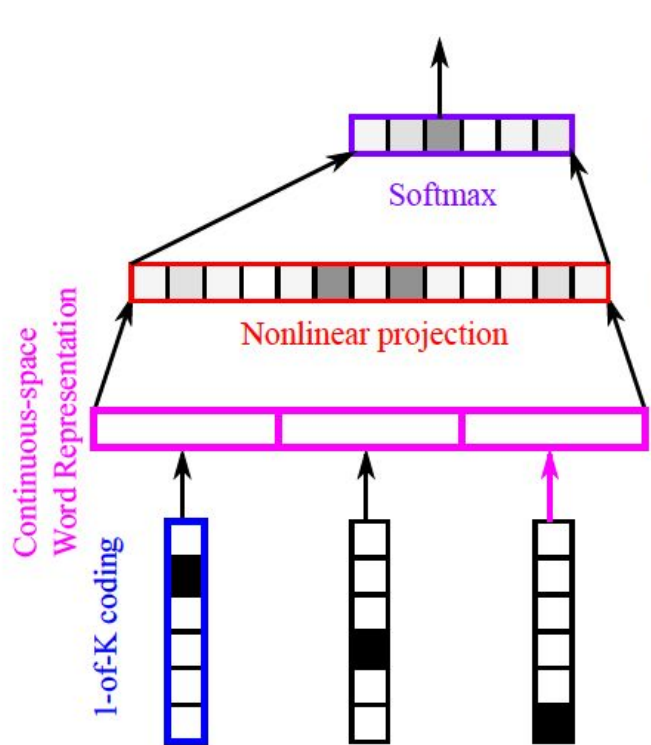
**Thanks ! Q&A ?**

<https://www.costa-jussa.com>





# Architecture: neural language model



$$p(w_n = k | w_1, w_2, \dots, w_{n-1}) = \mu_k = \frac{\exp(\mathbf{u}_k^\top \mathbf{h} + c_k)}{\sum_{k'=1}^K \exp(\mathbf{u}_{k'}^\top \mathbf{h} + c_{k'})}, \quad \mathbf{u}_k \in \mathbb{R}^{\dim(\mathbf{h})}$$

$$\sum_{k=1}^K \mu_k = 1.$$

$$\mathbf{h} = \tanh(\mathbf{W}\mathbf{p} + \mathbf{b})$$

$$\mathbf{E}^\top \mathbf{w}_i = \mathbf{e}_i. \quad \mathbf{E} = \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_{|V|} \end{bmatrix}, \quad \mathbf{p}^j = \mathbf{E}^\top \mathbf{w}^j \quad \mathbf{E} \in \mathbb{R}^{|V| \times d}.$$

$$\mathbf{w}_i = [0, 0, \dots, \underbrace{1}_{i\text{-th element}}, \dots, 0]^\top \in \{0, 1\}^{|V|}$$