

# Deep Learning for Natural Language Processing

## Lecture 5 — Text generation 1: Language models and word embeddings

---

Dr. Ivan Habernal

May 9, 2023

Trustworthy Human Language Technologies  
Department of Computer Science  
Technical University of Darmstadt



[www.trusthlt.org](http://www.trusthlt.org)

# The story so far

---

The story so far

- Recap 1: MLP and non-linearity

- Recap 2: Embedding categorical features

Language models

- Probability refresher

- 'Classical' language models

- Neural language models

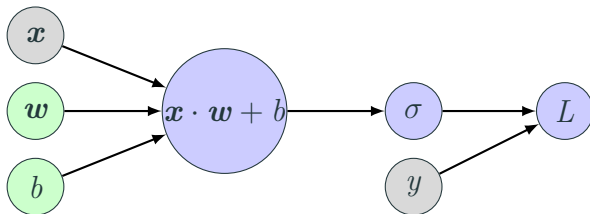
Word embeddings

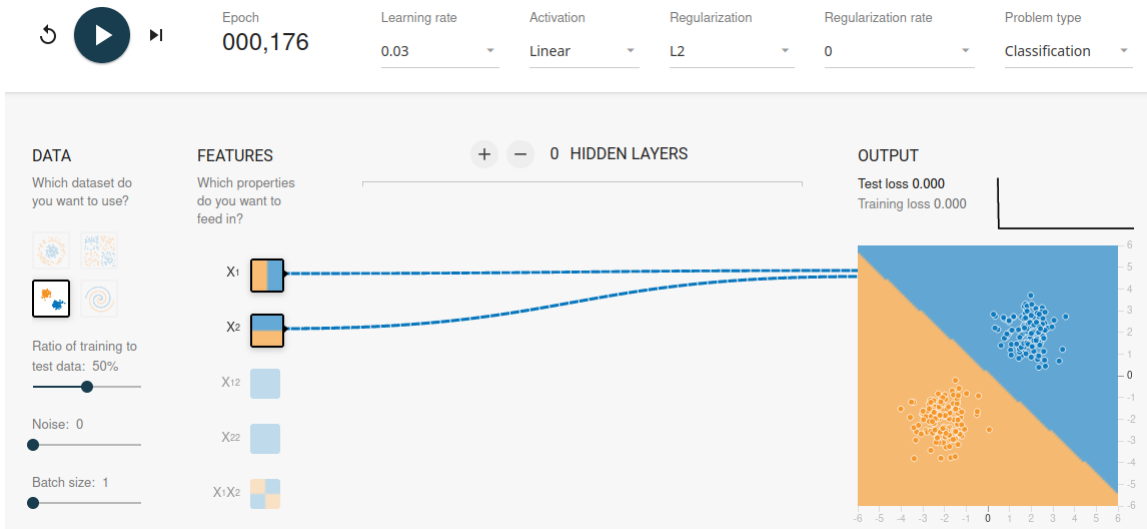
# The story so far

---

Recap 1: MLP and non-linearity

## Recap: Log-linear model for binary classification





**Figure 1:** Linear model can tackle only linearly-separable problems (<http://playground.tensorflow.org>)

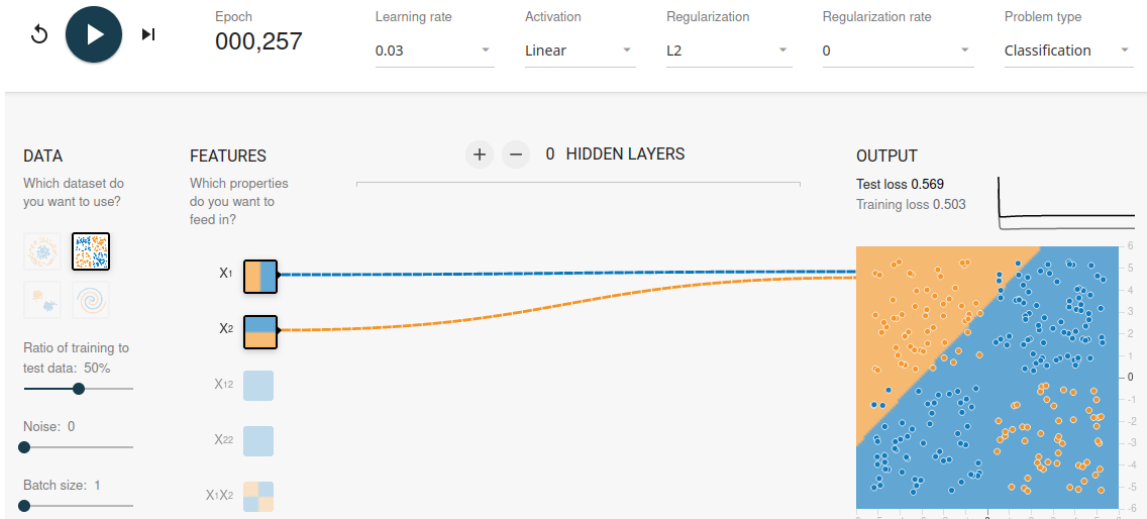
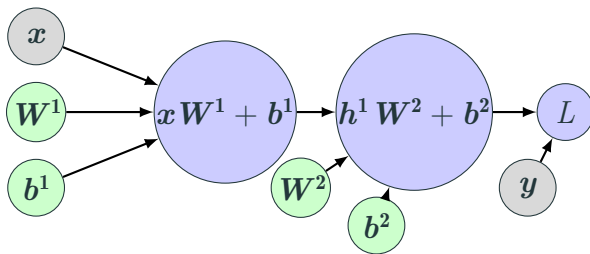


Figure 2: But will fail on, e.g., exclusive-or (XOR) tasks (<http://playground.tensorflow.org>)

## Recap: Stacking linear layers is still a linear model

$$\mathbf{x} \in \mathbb{R}^{d_{in}} \quad \mathbf{W}^1 \in \mathbb{R}^{d_{in} \times d_1} \quad \mathbf{b}^1 \in \mathbb{R}^{d_1} \quad \mathbf{W}^2 \in \mathbb{R}^{d_1 \times d_{out}} \quad \mathbf{b}^2 \in \mathbb{R}^{d_{out}}$$

$$f(\mathbf{x}) = (\mathbf{x}\mathbf{W}^1 + \mathbf{b}^1) \mathbf{W}^2 + \mathbf{b}^2$$



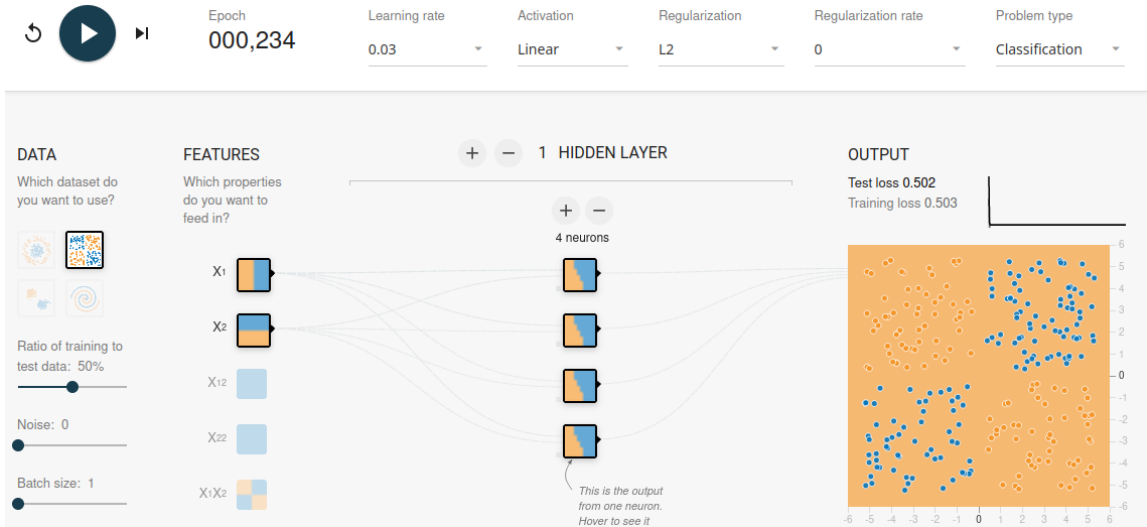
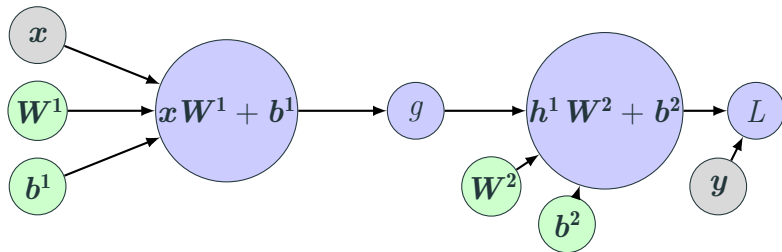


Figure 3: Linear hidden layers do not help  
(<http://playground.tensorflow.org>)



## Recap: Add non-linear function $g : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_1}$ (apply element-wise)

$$f(x) = g(xW^1 + b^1)W^2 + b^2$$



Typical non-linearities (activation functions),  $z \in \mathbb{R}$ :

- ReLU:  $\text{ReLU}(z) = \max(0, z)$
- tanh (hyperbolic tangent):  $\tanh(z) = \frac{e^{2z}-1}{e^{2z}+1}$

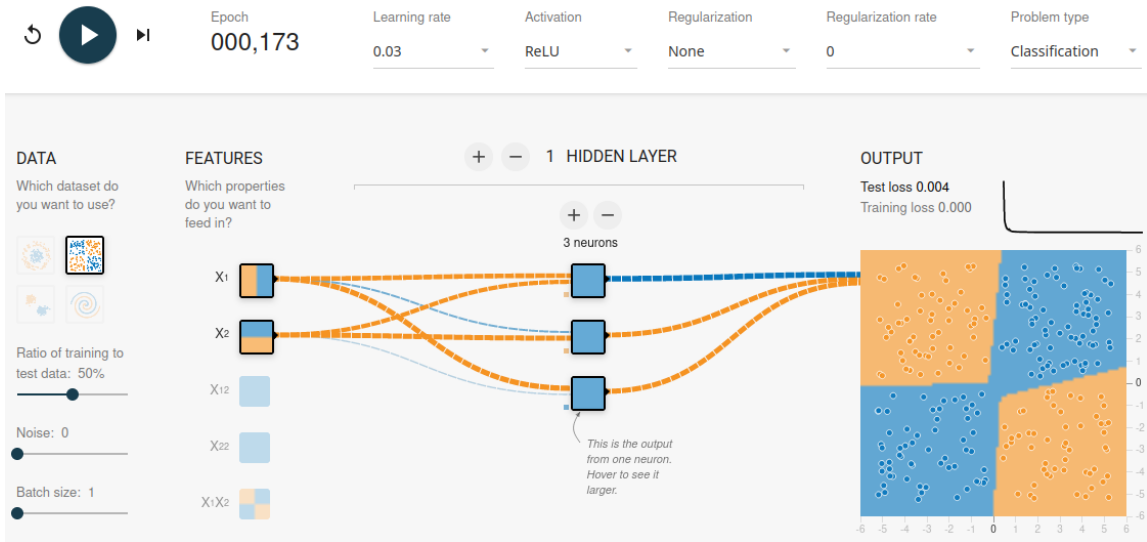


Figure 4: XOR solvable with, e.g., ReLU  
(<http://playground.tensorflow.org>)

# XOR example in super-simplified sentiment classification

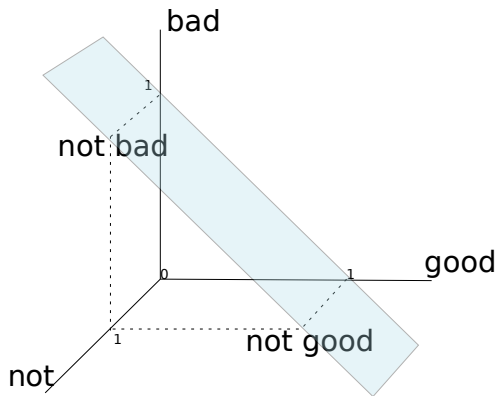


Figure 5:  $V = \{\text{not}, \text{bad}, \text{good}\}$ , binary features  $\in \{0, 1\}$

## The story so far

---

Recap 2: Embedding categorical features

## Recap: Matrix $W$ as learned representation — rows

$$\hat{y} = xW + b$$

	En	Fr	De	It	Es	Ot
a	•	•	•	•	•	•
...						
ZOO	•	•	•	•	•	•

Each of the  $d_{in}$  rows corresponds to a particular unigram, and provides a 6-dimensional vector representation

$$f(x) = g(xW^1 + b^1) W^2 + b^2$$

In MLP, the first layer's  $W^1$  learns representations ('embeddings') of the categorical features in  $x$

# Today: Language models and word embeddings

---

The story so far

- Recap 1: MLP and non-linearity

- Recap 2: Embedding categorical features

Language models

- Probability refresher

- 'Classical' language models

- Neural language models

Word embeddings

# Language models

---

The story so far

Recap 1: MLP and non-linearity

Recap 2: Embedding categorical features

Language models

Probability refresher

‘Classical’ language models

Neural language models

Word embeddings

# Language models

---

Probability refresher



# Probability refresher 1

## Categorical random variables

For example, the first word in a sentence

$W_1 \in \{\text{the, be, to, of, and, ...}\}$ , we assume a fixed vocabulary

## Probability distribution over random variables

For example, probability of '*the*' at position 1

$$\Pr(W_1 = w_1) = \Pr(W_1 = \text{the}) = 0.00024$$

Notation shortcuts:  $\Pr(W_1 = w_1) \rightarrow P(W_1), P(\text{the}), \text{etc.}$

# Probability refresher 2

## Joint probability

For example, probability of '*the*' at position 1 and '*cat*' at position 2

$$\Pr(W_1 = \text{the} \cap W_2 = \text{cat}) = 0.0000074$$

Notation shortcuts:  $P(W_1, W_2) = P(W_2, W_1)$

## Conditional probability

For example, probability of '*cat*' at position 2, **given** '*the*' at position 1

$$\Pr(W_2 = \text{cat} | W_1 = \text{the}) = \frac{P(W_1, W_2)}{P(W_1)}$$

# Probability refresher 3

## Independence

Two random variables  $X, Y$  are **independent** if and only if

$$P(X, Y) = P(X) \cdot P(Y)$$

## Conditional independence

Two random variables  $X, Y$  are **conditionally independent** given  $Z$  if and only if

$$P(X, Y|Z) = P(X|Z) \cdot P(Y|Z)$$

# Language models

---

‘Classical’ language models

# Goal of language modeling

Assign a probability to sentences in a language

## Example

“What is the probability of seeing the sentence *the lazy dog barked loudly*?”

Assigns a probability for the likelihood of given word (or a sequence of words) to follow a sequence of words

## Example

“What is the probability of seeing the word *barked* after the seeing sequence *the lazy dog*?”

# Language models formally

Sequence of words  $w_{1:n} = w_1 w_2 w_3 \dots w_n$  estimate

$$\Pr(w_{1:n}) = \Pr(w_1, w_2, \dots, w_n)$$

**Note:** we're sloppy in notation and usually omit the RVs

$$\Pr(W_1 = w_1, W_2 = w_2, \dots, W_n = w_n)$$

We *factorize* the joint probability into a product

- One factorization is very useful: left-to-right

$$\Pr(w_{1:n}) = \Pr(w_1 | \langle S \rangle) \Pr(w_2 | \langle S \rangle, w_1) \Pr(w_3 | \langle S \rangle, w_1, w_2) \dots \Pr(w_k | \langle S \rangle, w_1, w_2, \dots, w_{n-1})$$

# Simplifications in 'classical' language models

Despite factorization, the last term of  $\Pr(w_{1:n}) = \Pr(w_1|\langle s \rangle) \Pr(w_2|\langle s \rangle, w_1) \Pr(w_3|\langle s \rangle, w_1, w_2) \cdots \Pr(w_k|\langle s \rangle, w_1, w_2, \dots, w_{n-1})$  still depends on all the previous words of the sequence

## *k*-th order markov-assumption

The next word depends only on the last  $k$  words

$$\Pr(w_i|w_{1:i-1}) \approx \Pr(w_i|w_{i-k:i-1}) \quad (\text{inclusive indexing!})$$

$\langle s \rangle_0$  The<sub>1</sub> cat<sub>2</sub> sat<sub>3</sub> on<sub>4</sub> the<sub>5</sub>  $w_6$

$$i = 6, k = 2 \rightarrow \Pr(w_i|w_{1:i-1}) \approx \Pr(w_6|W_4 = \text{on}, W_5 = \text{the})$$

# Estimating probabilities in ‘classical’ language models

Maximum Likelihood Estimation (aka. counting and dividing)

$$\hat{P}_{\text{MLE}}(W_i = w | w_{i-k:i-1}) = \frac{\#(w_{i-k} \quad w_{i-k+1} \quad \dots \quad w_{i-1} \quad w)}{\#(w_{i-k} \quad w_{i-k+1} \quad \dots \quad w_{i-1})}$$

**What if  $\#(w_{i-k} \quad w_{i-k+1} \quad \dots \quad w_{i-1}) = 0$ ?**

ˆ Add-alpha smoothing ( $0 \leq \alpha \leq 1$ )

$$\hat{P}_{\text{add-}\alpha}(W_i = w | w_{i-k:i-1}) = \frac{\#(w_{i-k} \dots w_{i-1} \quad w) + \alpha}{\#(w_{i-k} \quad \dots \quad w_{i-1}) + \alpha |V|}$$



# Evaluating language models: Perplexity

Recall: Trained LM tells us probability of 'sentence'  $s$ :  $\Pr(s)$

Let's have  $n$  sentences in a test corpus, each of them has a uniform probability of appearing:  $\frac{1}{n}$

Then the **cross-entropy** (last lecture!) of our model is

$$\sum_{i=1}^n \frac{1}{n} \log\left(\frac{1}{\Pr(s_i)}\right) = \frac{1}{n} \sum_{i=1}^n \log\left(\frac{1}{\Pr(s_i)}\right) = -\frac{1}{n} \sum_{i=1}^n \log \Pr(s_i)$$

## Perplexity of LM

$$2^{\text{cross-entropy}} = 2^{(-\frac{1}{n} \sum_{i=1}^n \log \Pr(s_i))}$$

# Shortcomings of $n$ -gram language models

## Long-range dependencies

- To capture a dependency between the next word and the word 10 positions in the past, we need to see a relevant 11-gram in the text

Y. Goldberg (2017). *Neural Network Methods for Natural Language Processing*. Morgan & Claypool, p. 108

## Lack of generalization across contexts

- Having observed *black car* and *blue car* does not influence our estimates of the event *red car* if we haven't see it before

# Language models

---

Neural language models

Let's build a neural network

- Input: a  $k$ -gram of words  $w_{1:k}$
- Desired output: a probability distribution over the vocabulary  $V$  for the next word  $w_{k+1}$

## Embedding layer once again (recall last lecture)

If the input are symbolic **categorical features**

- e.g., words from a closed vocabulary

it is common to associate each possible feature value

- i.e., each word in the vocabulary

with a  $d$ -dimensional vector for some  $d$

These vectors are also *parameters* of the model, and are trained jointly with the other parameters

## Embedding layer: Lookup operation

The mapping from a symbolic feature values such as **word-number-48** to  $d$ -dimensional vectors is performed by an embedding layer (a lookup layer)

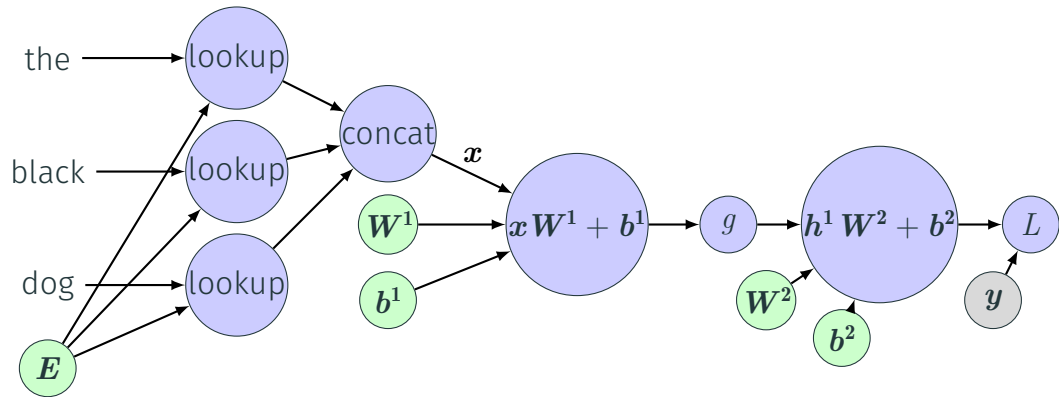
The parameters in an embedding layer are a matrix  $\mathbf{W}^{|V| \times d}$ , each row corresponds to a different word in the vocabulary

The lookup operation is then indexing  $v(w)$ , e.g.,

$$v(w) = v_{48} = \mathbf{E}_{[48,:]}$$

If the symbolic feature is encoded as a one-hot vector  $\mathbf{x}$ , the lookup operation can be implemented as the multiplication  $\mathbf{x}\mathbf{E}$

## Example network concatenating 3 words as embeddings ( $d_w = 50$ )



Each word  $\in \mathbb{R}^{|V|}$  (one hot),  $\mathbf{E} \in \mathbb{R}^{|V| \times 50}$ , each lookup output  $\in \mathbb{R}^{50}$ , concat output  $\mathbf{x} \in \mathbb{R}^{150}$

# Neural LMs

Let's build a neural network

- Input: a  $k$ -gram of words  $w_{1:k}$
- Desired output: a probability distribution over the vocabulary  $V$  for the next word  $w_{k+1}$

Each input word  $w_k$  is associated with an embedding vector  $v(w) \in \mathbb{R}^{d_w}$  ( $d_w$  — word embedding dimensionality)

Input vector  $\mathbf{x}$  is a concatenation of  $k$  words

$$\mathbf{x} = [v(w_1); v(w_2); \dots; v(w_k)]$$



MLP with one (or more) hidden layers

$$v(w) = \mathbf{E}_{w,:}$$

$$\mathbf{x} = [v(w_1); v(w_2); \dots; v(w_k)]$$

$$\mathbf{h} = g(\mathbf{x} \mathbf{W}^1 + \mathbf{b}^1)$$

$$\hat{\mathbf{y}} = \Pr(W_i | w_{1:k}) = \text{softmax}(\mathbf{h} \mathbf{W}^2 + \mathbf{b}^2)$$

Output dimension:  $\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$

# Training neural LMs

Where to get training examples?

Training examples are simply word  $k$ -grams from an unlabeled corpus

- Identities of the first  $k - 1$  words are used as features
- The last word is used as the target label for the classification

The model is trained using cross-entropy loss

## Some advantages and limitations of neural LMs

$\approx$  linear increase in parameters with  $k + 1$  (better than 'classical' LMs) but

- The size of the output vocabulary affects the computation time
- The softmax at the output layer requires an expensive matrix-vector multiplication with the matrix  $\mathbf{W}^2 \in \mathbb{R}^{d_{\text{hid}} \times |V|}$ , followed by  $|V|$  exponentiations

Solutions: Hierarchical softmax, noise-contrastive estimation

# Generating text with language models

We can generate (“sample”) random sentences from the model according to their probability

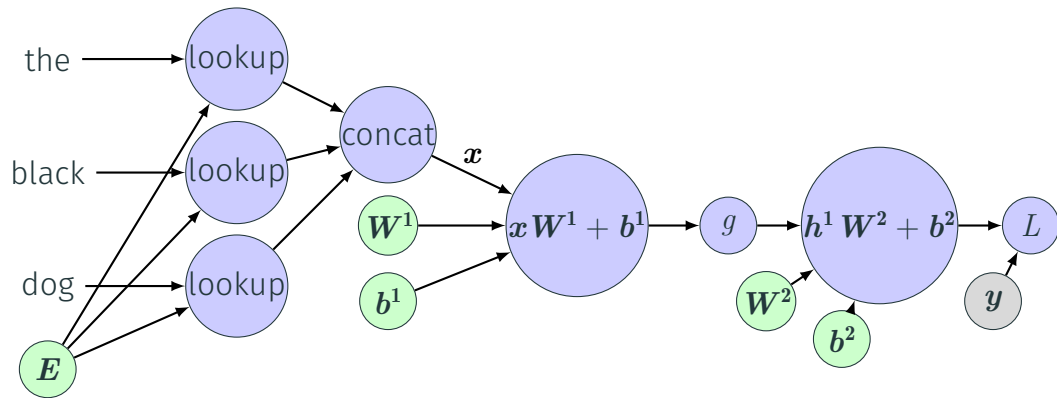
1. Predict a probability distribution over the vocabulary conditioned on the start symbol <s>
2. Draw a random word (the first word) according to the predicted distribution
3. Predict a probability distribution over the vocabulary conditioned on the start symbol and the first word
4. Draw a random word (the second word) according to the predicted distribution
5. Repeat until generated *end-of-sentence* symbol </s> (or <EOS>)

# Sampling words — alternatives

Sampling (generating) the most probable word at each step might not be optimal globally

- Beam search — generate top  $k$  candidates at each step

# Learned word representations as a by-product



Each row of  $E$  learns a word representation

Each column of  $W^2$  learns a word representation

# Word embeddings

---

The story so far

- Recap 1: MLP and non-linearity

- Recap 2: Embedding categorical features

Language models

- Probability refresher

- 'Classical' language models

- Neural language models

Word embeddings

# Word embeddings as pre-trained word representation

Option A: We can initialize the embeddings matrix  $\mathbf{E}$  randomly and learn during our supervised task

Option B: Use pre-trained word embeddings from task for which we have a lot of data

- Use self-supervised learning (create labeled data ‘for free’ using the next word prediction objective)
- Learned word embedding matrix plugged into our supervised task

Desired word embeddings properties: ‘Similar’ words have similar embeddings vectors



# Recap

---

The story so far

- Recap 1: MLP and non-linearity

- Recap 2: Embedding categorical features

Language models

- Probability refresher

- 'Classical' language models

- Neural language models

Word embeddings

# Take aways

- Language modeling is an essential part of contemporary NLP
- Self-supervised models, unlabeled data, next word prediction
- Neural language models learn embedding of words

# License and credits

Licensed under Creative Commons  
Attribution-ShareAlike 4.0 International  
(CC BY-SA 4.0)



## Credits

Ivan Habernal

Content from ACL Anthology papers licensed under CC-BY

<https://www.aclweb.org/anthology>

XOR examples generated by <http://playground.tensorflow.org/>