

# Language Models

Joost Bastings

<http://joost.ninja>

you know nothing, jon

life is like a box of

the

the quick

the quick brown

the quick brown fox

the quick brown fox jumps



the quick brown fox jumps over

the quick brown fox jumps over the

the quick brown fox jumps over the lazy

the quick brown fox jumps over the lazy dog

# Predicting the future

Language models give us the probability of a **sentence**.  
At any time step, they assign a **probability** to the next word.

# Applications

This is useful when we deal with **noisy input**:

- speech recognition
- handwriting recognition
- spelling correction
- machine translation

**today:**

n-gram based Im

feature based Im

neural network based Im

## Probability of a sentence

$$p(\text{you know nothing jon snow}) = ???$$

We can use the **chain rule**:

$$p(B | A) = p(A, B) / p(A)$$

$$p(\text{you know nothing jon snow}) = \\ p(\text{know} | \text{you}) p(\text{nothing} | \text{you know}) p(\text{jon} | \text{you know nothing}) p(\text{snow} | \text{you know nothing jon})$$



## n-gram language models

To make (count based) estimation possible, we need a **markov assumption**

**bi-gram** case:

$$p(x_i | x_1 x_2 \dots x_{i-1}) \approx p(x_i | x_{i-1})$$

parameters:  $\theta_{x|x'} \quad \forall x, x' \in \mathcal{V}$

**n-gram** case:

$$p(x_i | x_1 x_2 \dots x_{i-1}) \approx p(x_i | x_{i-n+1} \dots x_{i-1})$$

parameters:  $\theta_{x|h} \quad \forall x \in \mathcal{V}, h \in \mathcal{V}^{n-1}$

# Smoothing

If we didn't observe a certain bigram, then  $p(x_i | x_{i-1})$  will be 0

This makes the probability of the sentence also 0!

**MLE:**  $p_{\text{MLE}}(x_i | x_{i-1}) = \text{count}(x_{i-1}, x_i) / \text{count}(x_{i-1})$

**Laplace / add-one smoothing:**  $p_{\text{Add1}}(x_i | x_{i-1}) = \frac{\text{count}(x_{i-1}, x_i) + 1}{\text{count}(x_{i-1}) + V}$

*This doesn't work so well for language modeling.*

*For more advanced smoothing, cf. Kneser-Ney, Stupid Backoff.*

We would like our models to assign high probability to “real” sentences

Let  $\mathbf{x} = x_1 x_2 \dots x_N$

$$pp(\mathbf{x}) = \exp^{-1/N \sum_i \ln P(x_i)}$$

$$pp(\mathbf{x}) = P(x_1 x_2 \dots x_N)^{-1/N}$$

## Log-linear models

$$p_w(Y=y \mid X=x) = \frac{\exp \mathbf{w} \cdot \phi(x, y)}{\sum_{y' \in \mathcal{Y}} \exp \mathbf{w} \cdot \phi(x, y')}$$

Y is the **next word** (and  $\mathcal{Y}$  the vocabulary)

X is the **history**

$\phi$  is a **feature function** which returns a d-dimensional vector

$\mathbf{w}$  are the model parameters

## Why use log-linear models with features?

With *features* of words and histories we can share **statistical weight**

*Mind: with n-grams, there is no sharing at all!*

We also get **smoothing**

We can add arbitrary features, e.g. we could condition the text on properties of the author (gender, age, political affiliation,...)

We can easily **interpret** these kinds of models

Each feature  $\phi_k$  controls a *factor* to the probability ( $e^{w_k}$ )

if  $w_k < 0$ , then  $\phi_k$  makes the event *less likely* by a factor of  $1/e^{w_k}$

if  $w_k > 0$ , then  $\phi_k$  makes the event *more likely* by a factor of  $e^{w_k}$

if  $w_k = 0$ , then  $\phi_k$  has no effect

## What features should we use?

n-gram features: " $X_{j-1} = \text{the} \wedge X_j = \text{puppy}$ "

gappy n-gram features: " $X_{j-2} = \text{the} \wedge X_j = \text{puppy}$ "

spelling features: " $X_j$ 's first character is Capitalized"

class features: " $X_j$  is a member of class 123"

gazetteer features: " $X_j$  is a geographic place name"

# Comparison: n-gram vs. log-linear

## n-gram

$$p_{\theta}(\mathbf{x}) = \prod_i \theta_{x_i | \mathbf{h}_i}$$

parameters:

$$\theta_{x|h} \quad \forall x \in \mathcal{V}, \mathbf{h} \in \mathcal{V}^{n-1}$$

**MLE:**

$$\theta_{x|h}^{\wedge} = \text{count}(\mathbf{h}, x) / \text{count}(\mathbf{h})$$

## log-linear

$$p_{\theta}(\mathbf{x}) = \prod_i \exp \mathbf{w} \cdot \phi(\mathbf{h}_i, x_i) / Z_{\mathbf{w}}(\mathbf{h}_i)$$

parameters:

$$\mathbf{w}_k \quad \forall k \in \{1, \dots, d\}$$

**MLE:**

**No closed form! → use SGD**

training data  $\{(\mathbf{x}_i, \mathbf{h}_i)\}_{i=1}^N$

$$\text{MLE: } \max_{\mathbf{w}} \sum_i \mathbf{w} \cdot \phi(\mathbf{h}_i, x_i) - \log Z_{\mathbf{w}}(\mathbf{h}_i)$$

$$\nabla_{\mathbf{w}} f_i = \phi(\mathbf{h}_i, x_i) - \sum_v p(v | \mathbf{h}_i) \cdot \phi(\mathbf{h}_i, v)$$

Observed features

Expected features

# Stochastic Gradient Descent

**Goal:** minimize  $\sum_i^N f_i(\mathbf{w})$  with respect to  $\mathbf{w}$

**Input:** initial values  $\mathbf{w}$ , learning rate  $\alpha$

**SGD:**

**while** “not converged”:

sample training instance  $i$  from  $\{1, \dots, N\}$

$$\mathbf{w} \leftarrow \mathbf{w} - \alpha \cdot \nabla_{\mathbf{w}} f_i$$



# Regularization

## Problem:

if a certain feature is often positive, we can increase the objective by making its weight slightly larger every time

This can cause **overfitting**!

## Solution:

Regularization

We add a term  $\lambda ||\mathbf{w}||^2$  to our objective (“L2-norm”) to discourage large weights

# Summary so far

## n-gram

$\mathbf{h}_i$  is n-1 previous words

estimation: count and normalize

think about smoothing

## log-linear

featurized representation of  $(\mathbf{h}_i, x_i)$

estimation with SGD

think about features

# Neural language models

# Motivation

n-gram language models have proven to be effective for various tasks 🙌

*maybe we do not need the **full** history*

log-linear models allow sharing statistical weight through features 🙌

*sharing is important for generalization*

maybe our history is still too limited ( $n-1$  words) 🙌

*sometimes larger histories are useful, but hard to estimate parameters*

we need to find useful features 🙌

*and be careful of overfitting with too many features*

## Motivation (2)

Consider:

1. The cat is walking in the bedroom
2. A dog was running in a room

Observing (1) should help us generalize to make (2) almost as likely.

But with the previous methods this is hard to accomplish!

## Feed-forward NN approach

with neural networks we can exploit **distributed representations** to allow statistical weight sharing

*how it works:*

1. each word is assigned a distributed m-dimensional **feature vector**
2. a **probability function** of word sequences is expressed in terms of those vectors
3. we **jointly** learn the feature vectors and the parameters of that probability function

## Why would this work?

Similar words are expected to have a similar feature vector

*(dog, cat), (is, was), (running, walking), (room, bedroom), ...*

With this, probability mass is naturally transferred from (1) to (2):

1. *The cat is walking in the bedroom*

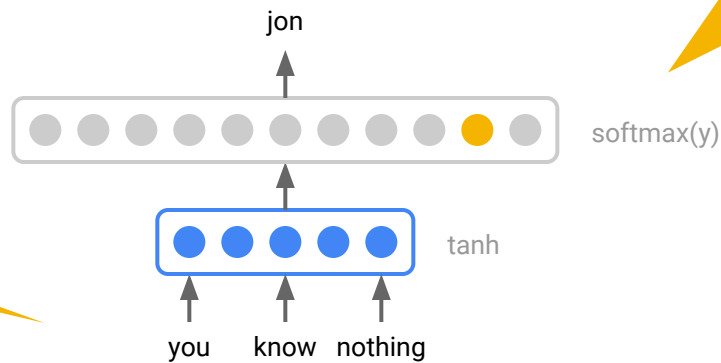
2. *A dog was running in a room*

and to many, many other similar sentences.

**The point:** the presence of only 1 sentence in the training data will increase the probability for a *combinatorial number* of “neighbors” in sentence space.

# Feed-forward NN language model

This is still an n-gram model!



Normalization over the full vocabulary.

$$y = W'' \tanh(Wx + b) + W'x + b'$$



## Why does it work?

the nonlinear NN function allows **feature combinations** a linear model cannot get

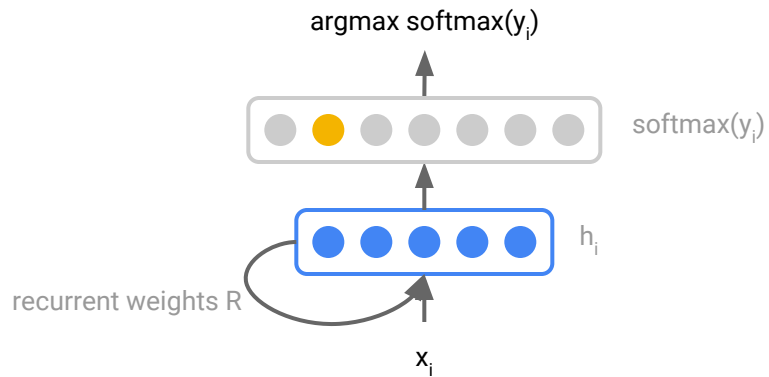
**end-to-end training** on next word prediction

## Recurrent NN language model

we now have much better generalization, but still a limited context

**recurrent neural networks** (RNNs) have an **unlimited context**

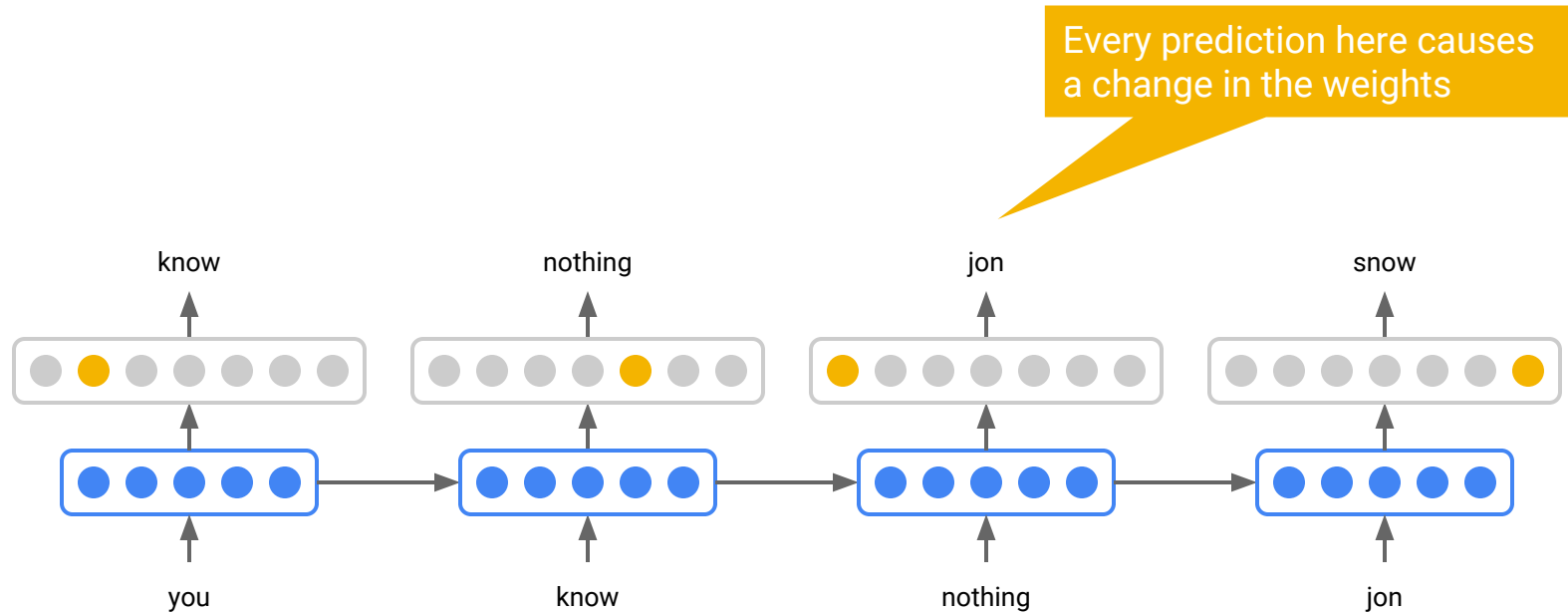
# Recurrent NN language model



$$h_i = \sigma( Wx + Rh_{i-1} + b )$$

$$y_i = W'' h_i + b'$$

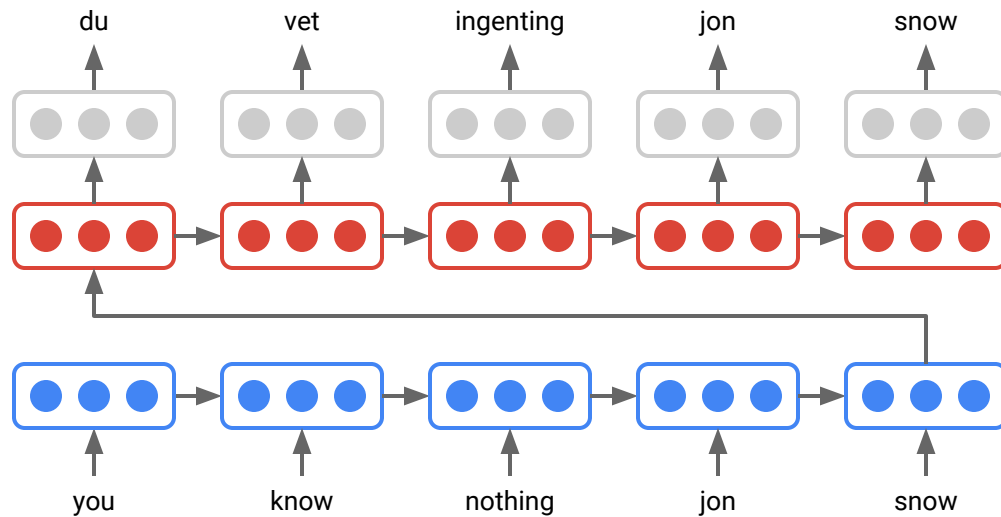
# Recurrent NN language model



## Final note on neural language models

1. the log-likelihood function is **not concave**
2. so when we evaluate a new neural language **model**,  
we also evaluate our **algorithm** to estimate the parameters  
(and its **hyperparameters**)
3. RNNs suffer from the **vanishing gradient problem** (next lecture)
4. many improved language models have been proposed since these ones  
e.g. log-bilinear lm, LSTM-based, character based

## Preview: Encoder-decoder



## References

**Bengio**, Yoshua, et al. "A neural probabilistic language model." Journal of machine learning research 3.Feb (2003): 1137-1155.

**Mikolov**, Tomas, et al. "Recurrent neural network based language model." Interspeech. Vol. 2. 2010.

**Cho**, Kyunghyun, et al. "Learning phrase representations using RNN encoder-decoder for statistical machine translation." arXiv preprint arXiv:1406.1078 (2014).

Noah Smith, CSEP517:

<http://courses.cs.washington.edu/courses/csep517/17sp/>

Jurafsky & Martin. Speech and Language Processing (3rd ed. draft).

<https://web.stanford.edu/~jurafsky/slp3/>