

Hidden Markov Models

Sequence labeling

Problem

- Given a sequence of *tokens*, infer the most probable sequence of *labels* for these tokens.

Examples

- part of speech tagging
- named entity recognition

Example: part-of-speech tagging

Tokens are words, and labels are PoS tags.

PRON

VERB

DET

PROPN

NOUN

I

saw

a

Heffalump

today.

PoS tags from Universal Dependencies project

1.0000000000

69: 0.8178592142
70: 0.8045185384

Open class words		Closed class words	
ADJ	adjective	ADP	adposition
ADV	adverb	AUX	auxiliary verb
INTJ	interjection	CCONJ	coordinating conjunction
NOUN	noun	DET	determiner
PROPN	proper noun	NUM	numeral
VERB	verb	PART	particle
Other		PRON	
PUNCT	punctuation	SCONJ	subordinating conjunction
SYM	symbol		
X	other		

<http://universaldependencies.org/>

Example: named entity recognition

Once upon a time, a very long time ago now, about

[last Friday], **[Winnie-the-Pooh]** lived in a forest all by
himself under the name of **[Sanders]**.

Types of named entities

Any real-world object which may have a proper name:

- persons
- organizations
- locations
- ...

Also named entities usually include:

- dates and times
- units
- amounts

Approaches to sequence labeling

1. Rule-based models (example: EngCG tagger)
2. Separate label classifiers for each token
3. Sequence models (HMM, MEMM, CRF)
4. Neural networks

PoS tagging with HMMs

Notation:

$\mathbf{x} = x_1, \dots, x_T$ is a sequence of words (input)

$\mathbf{y} = y_1, \dots, y_T$ is a sequence of their tags (labels)

We need to find the most probable sequence of tags given the sentence:

$$\mathbf{y} = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{x}, \mathbf{y})$$

But first, let us define the model...

Hidden Markov Model

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y}) p(\mathbf{y}) = \prod_{t=1}^T p(x_t|y_t) p(y_t|y_{t-1})$$

observable hidden

Markov assumption:

$$p(\mathbf{y}) \approx \prod_{t=1}^T p(y_t|y_{t-1})$$

Output independence:

$$p(\mathbf{x}|\mathbf{y}) \approx \prod_{t=1}^T p(x_t|y_t)$$

Text generation in HMM

Assume that the text is generated in the following manner:

- One chooses the next PoS tag given the previous tag
- Given the current tag, one generates another word

Thus, the neighboring words do not depend on each other, but they depend on the underlying tags.

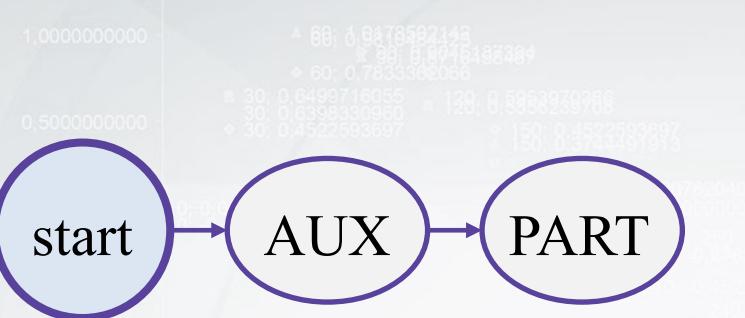
Text generation in HMM: an example



Text generation in HMM: an example

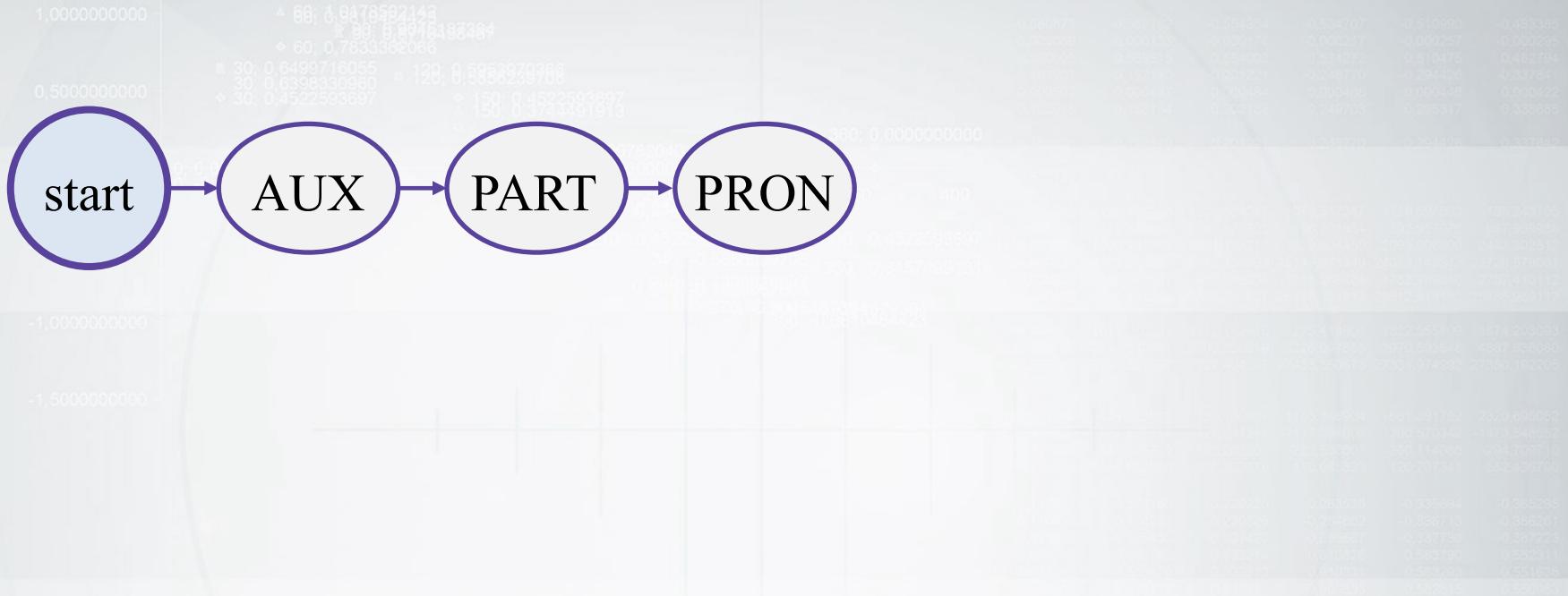


Text generation in HMM: an example

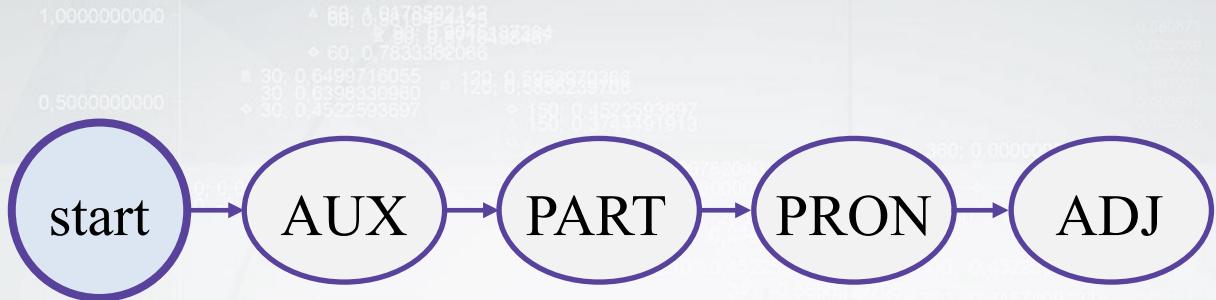


360; 0.0000000000	0.229220	0.263536	-0.335694	0.385298
360; 0.0000000000	-0.230325	-0.264602	-0.336713	0.386261
360; 0.0000000000	-0.231423	0.265687	-0.337730	0.387223
360; 0.0000000000	0.232524	0.261052	0.533780	0.352311
360; 0.0000000000	0.233195	0.263381	0.610331	0.533203
360; 0.0000000000	0.234200	0.262410	0.609533	0.532515
360; 0.0000000000	0.235174	0.263241	0.609583	0.530983

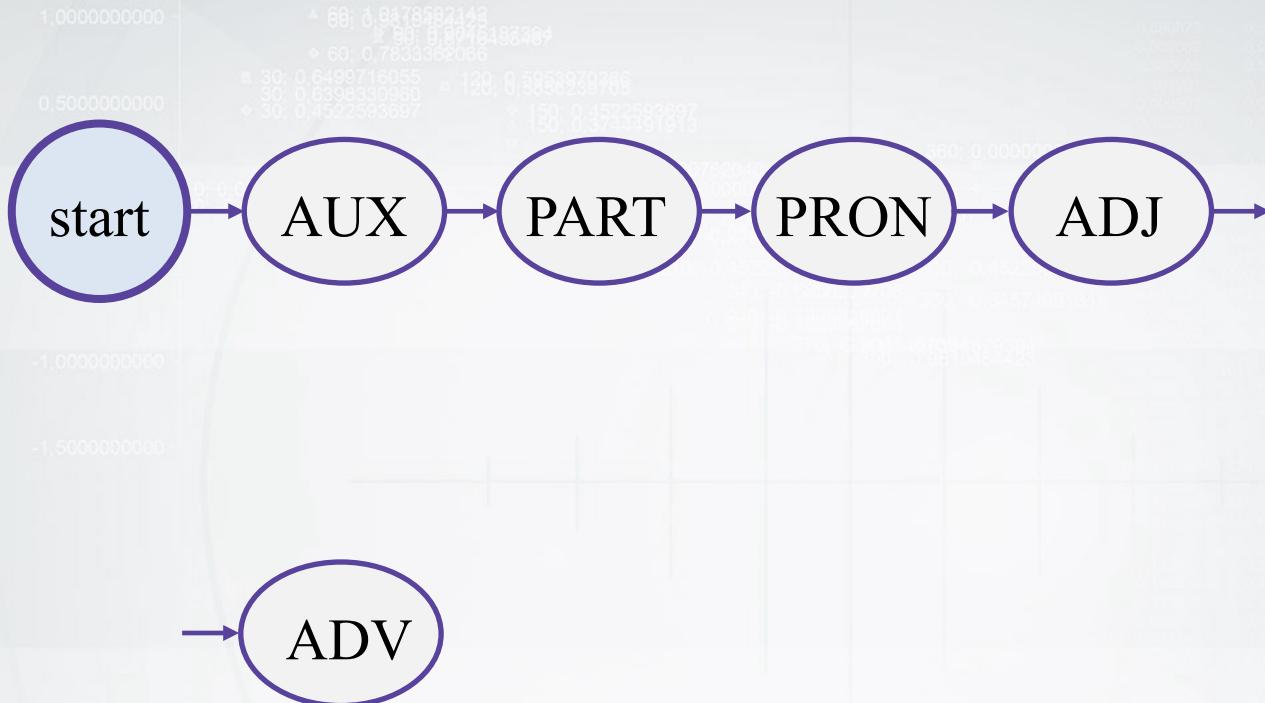
Text generation in HMM: an example



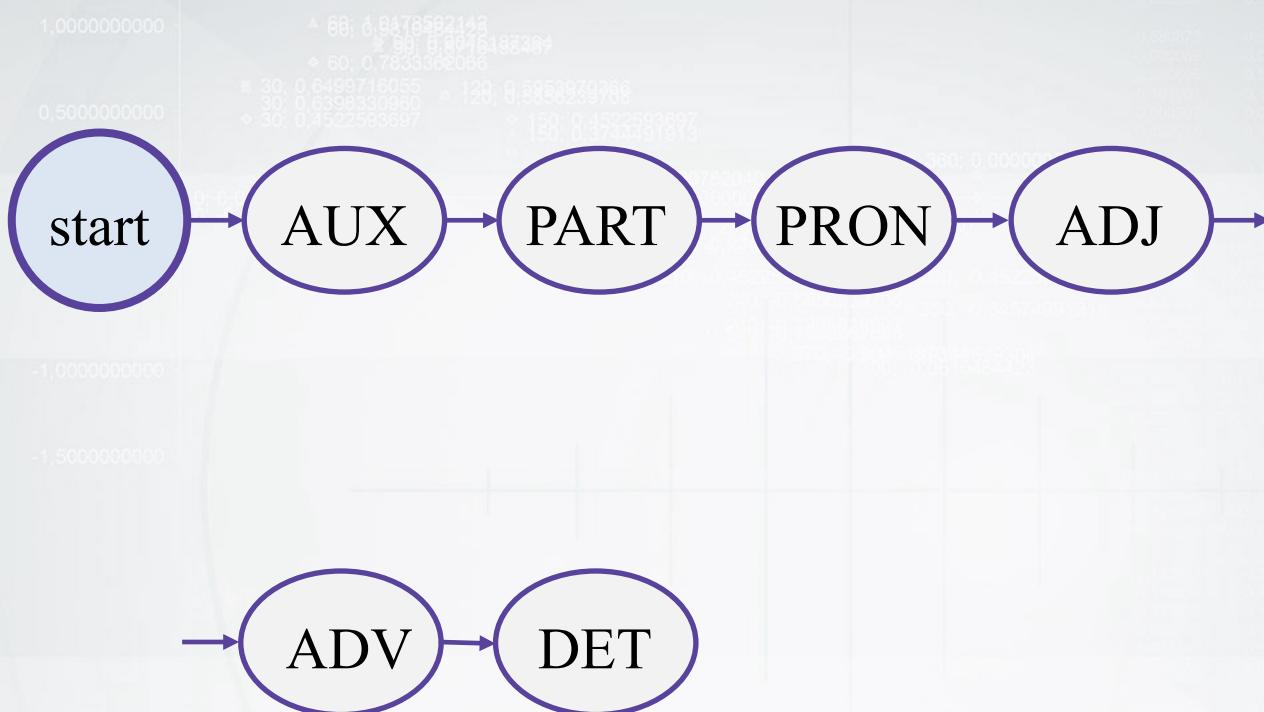
Text generation in HMM: an example



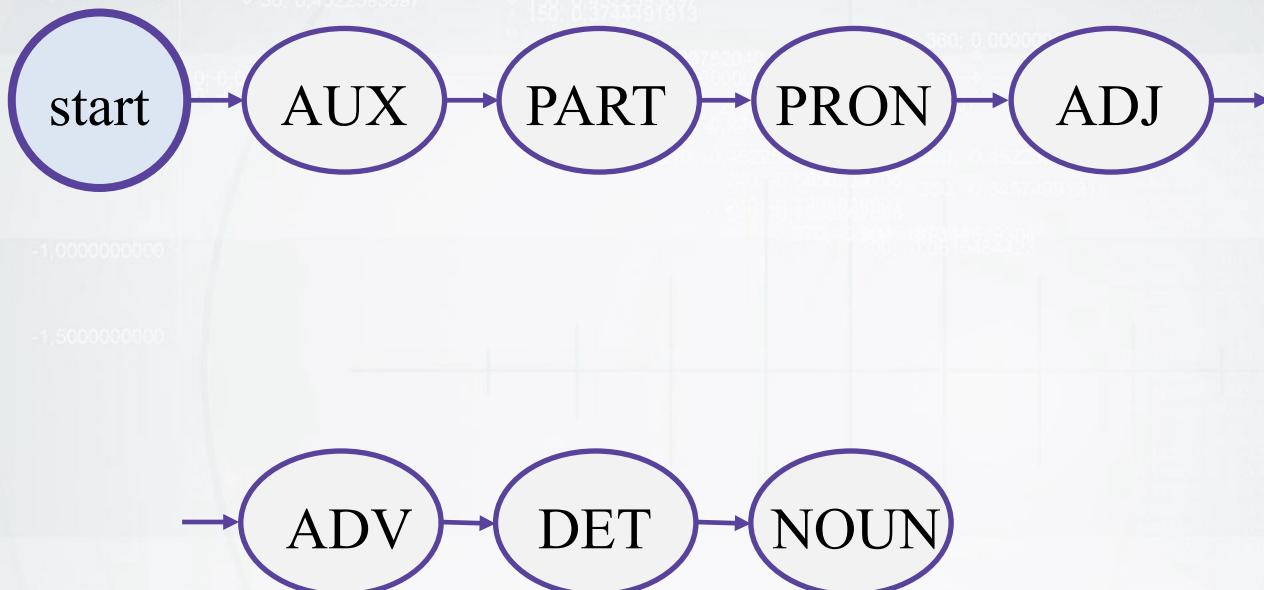
Text generation in HMM: an example



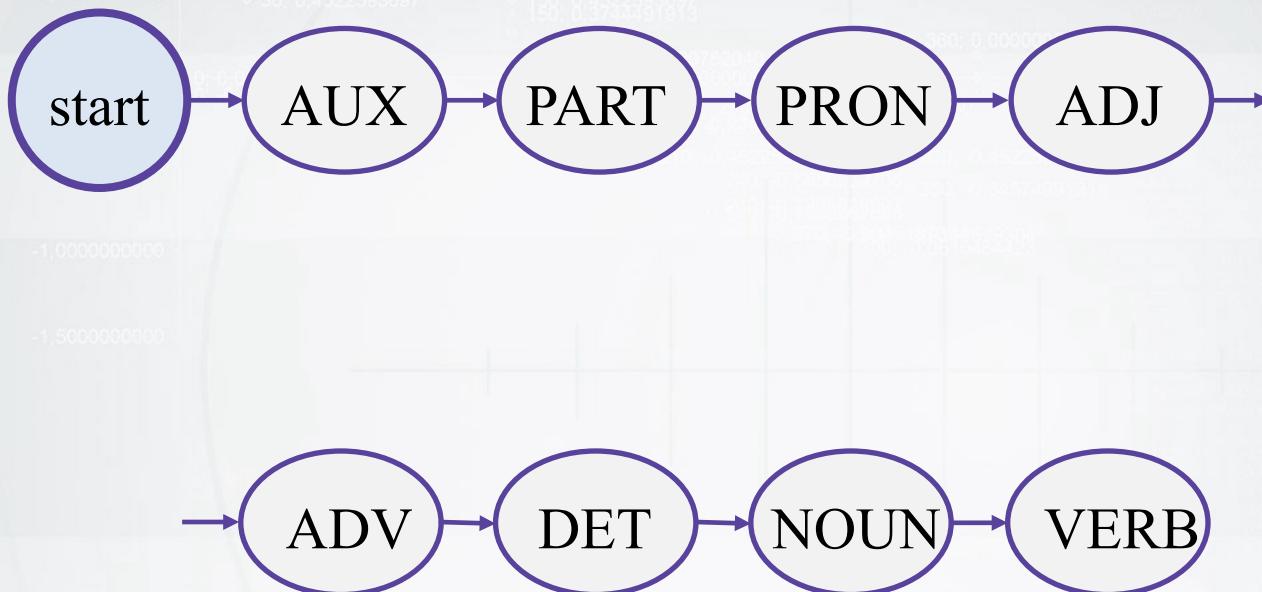
Text generation in HMM: an example



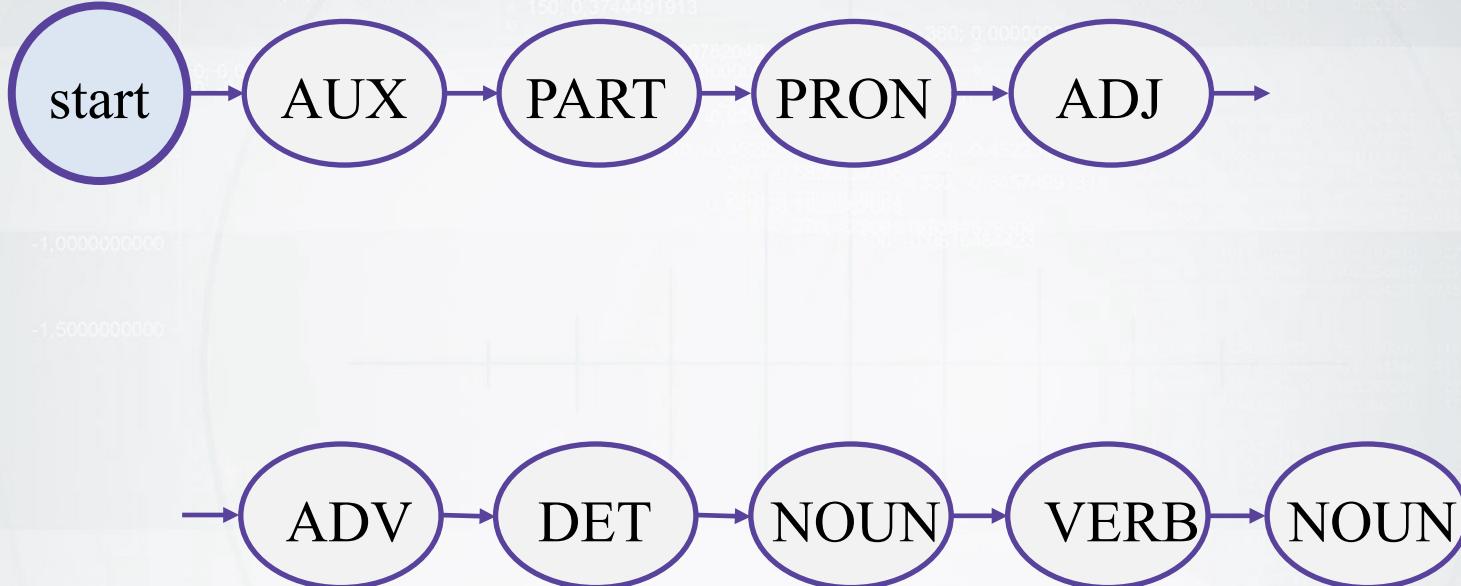
Text generation in HMM: an example



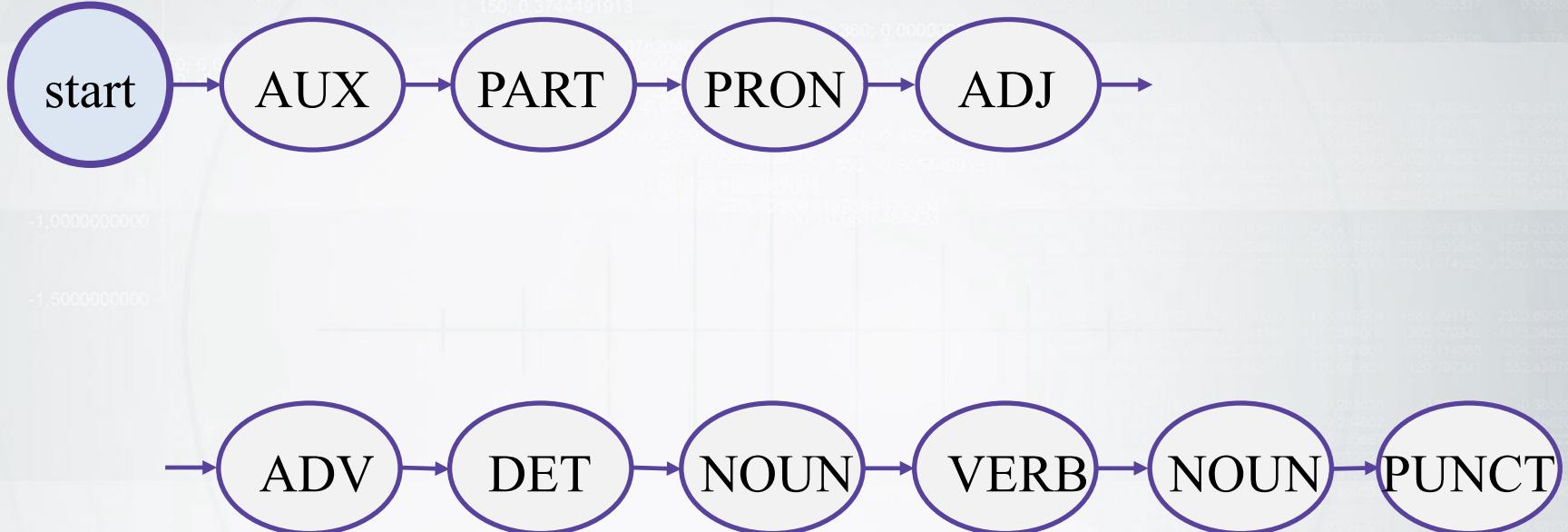
Text generation in HMM: an example



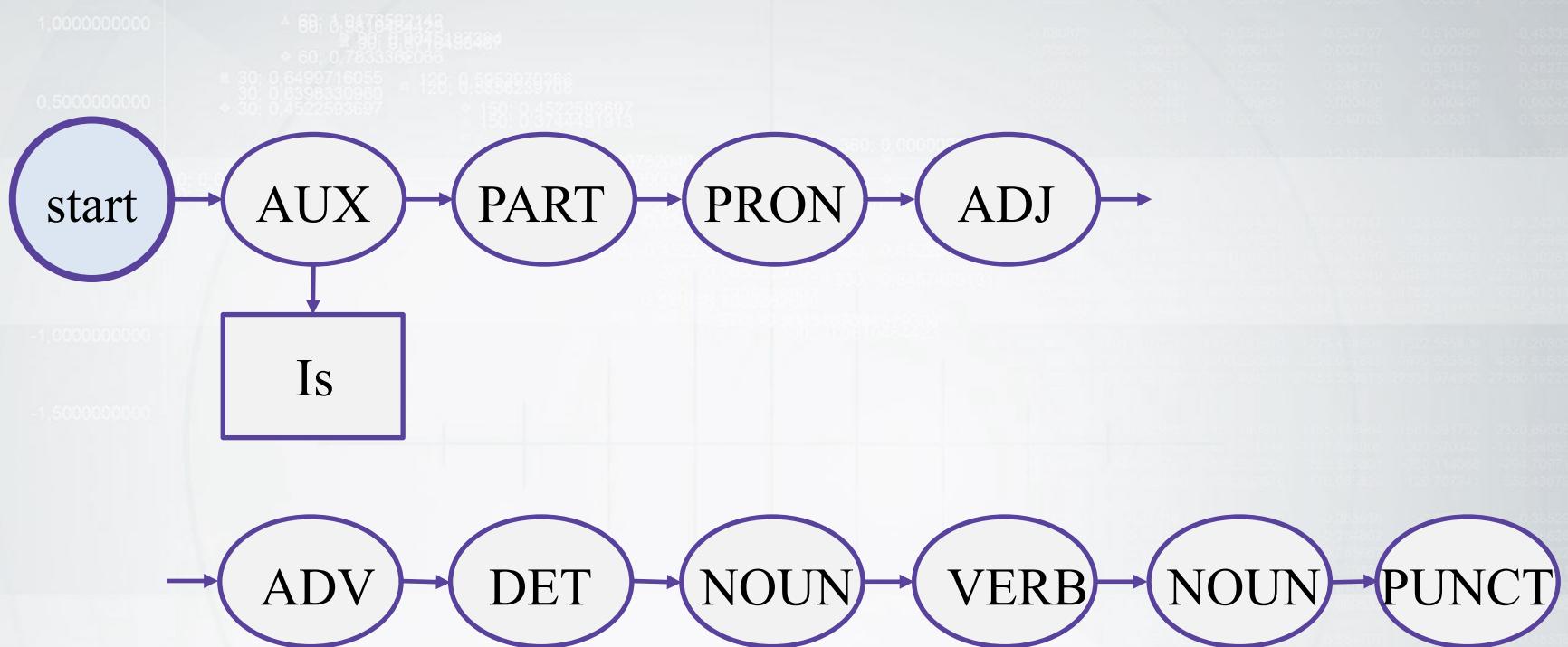
Text generation in HMM: an example



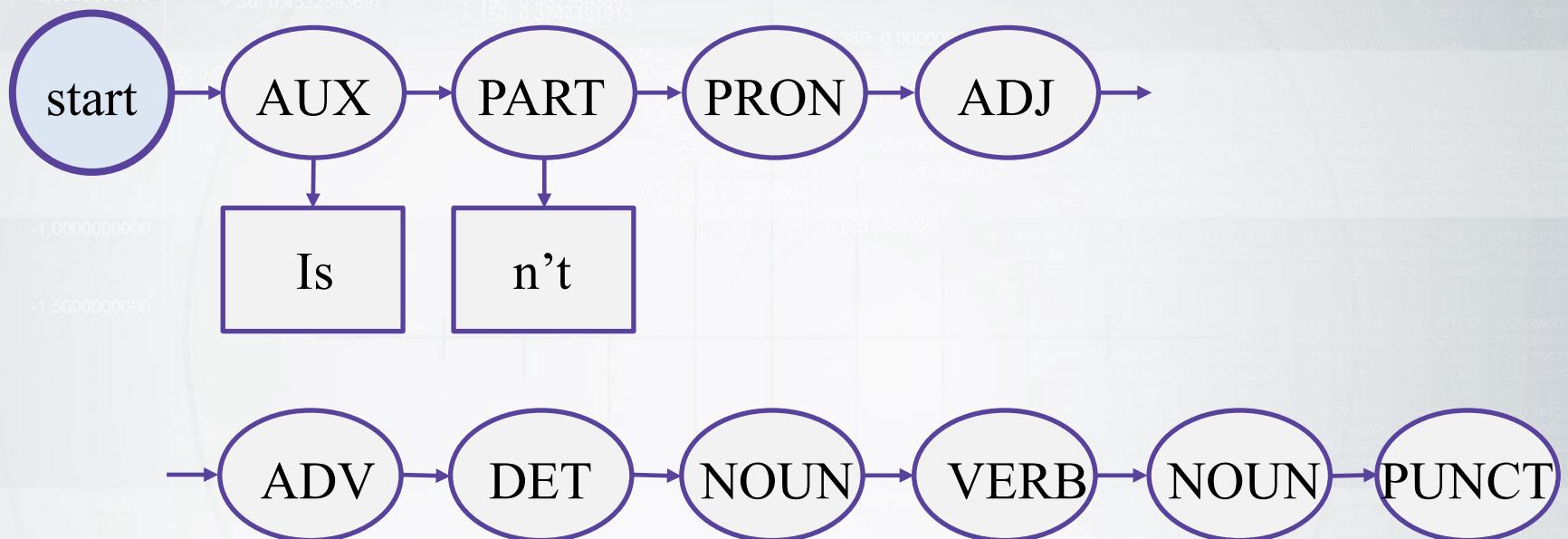
Text generation in HMM: an example



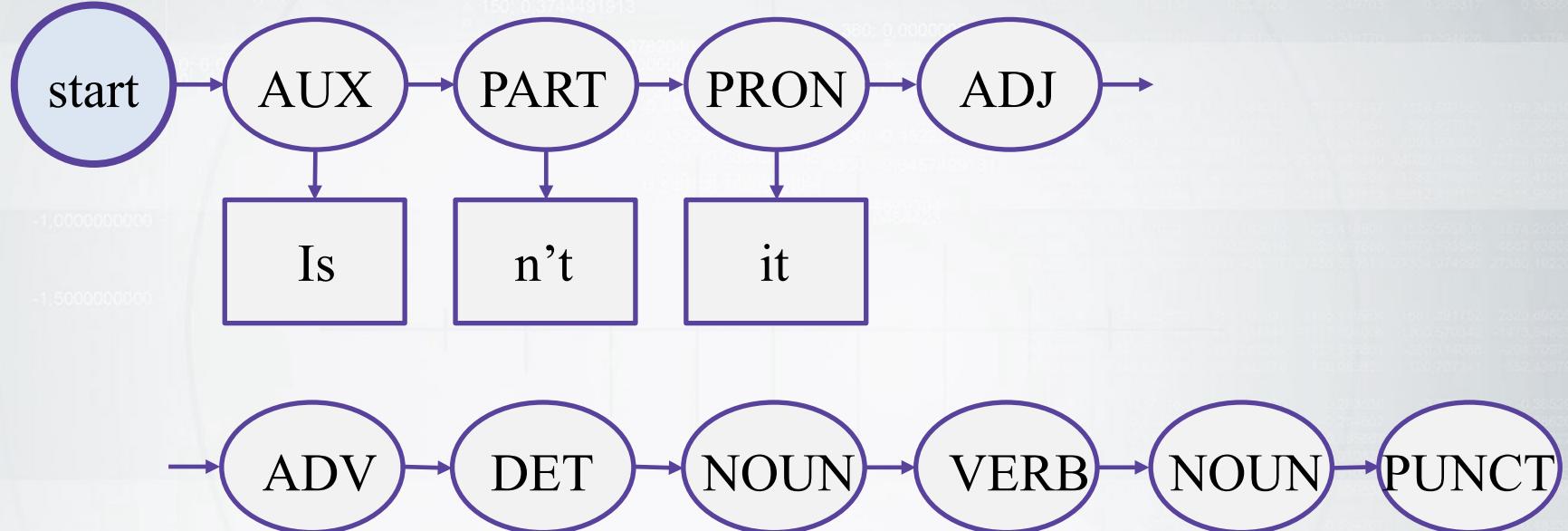
Text generation in HMM: an example



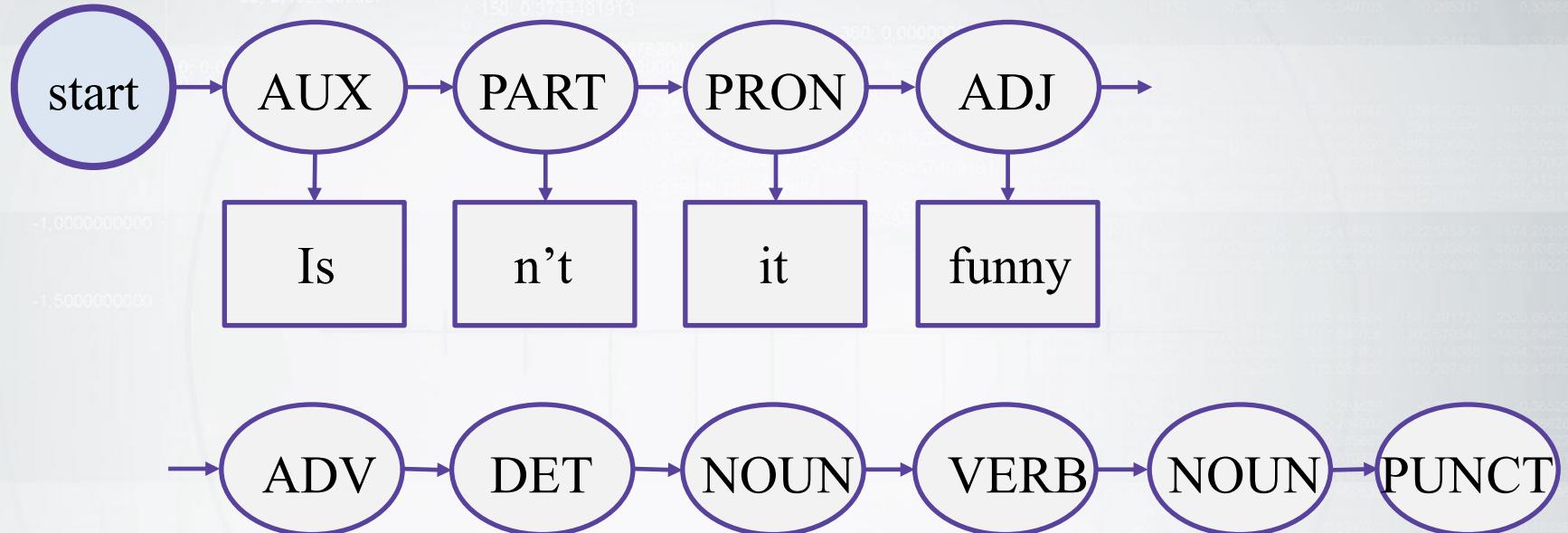
Text generation in HMM: an example



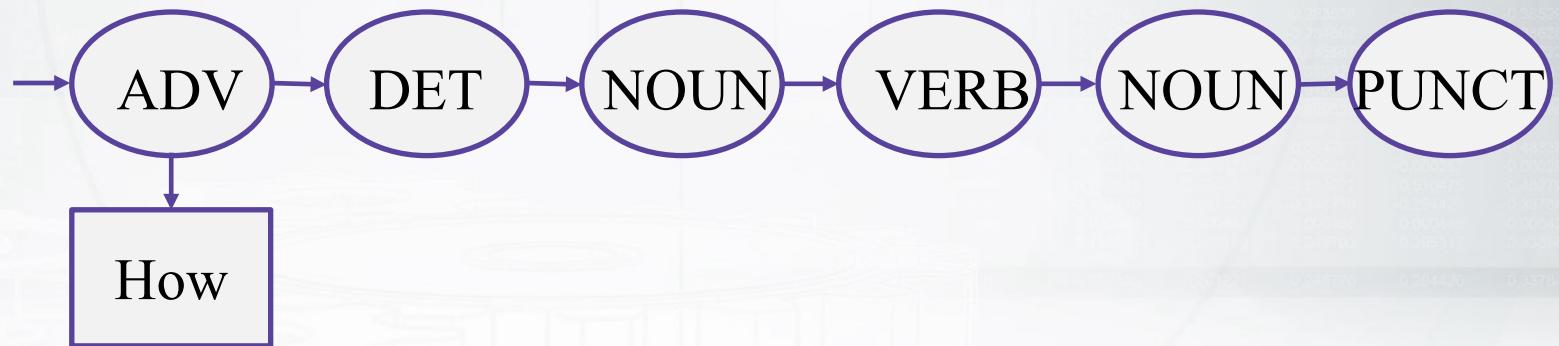
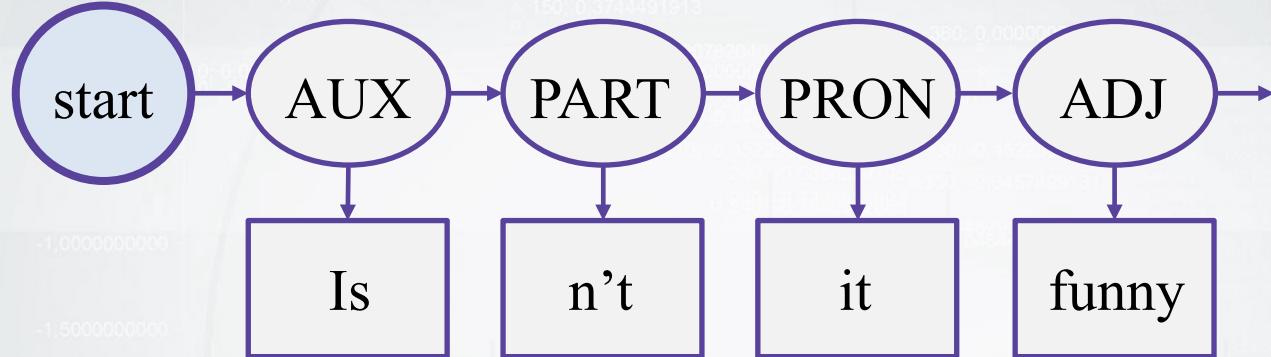
Text generation in HMM: an example



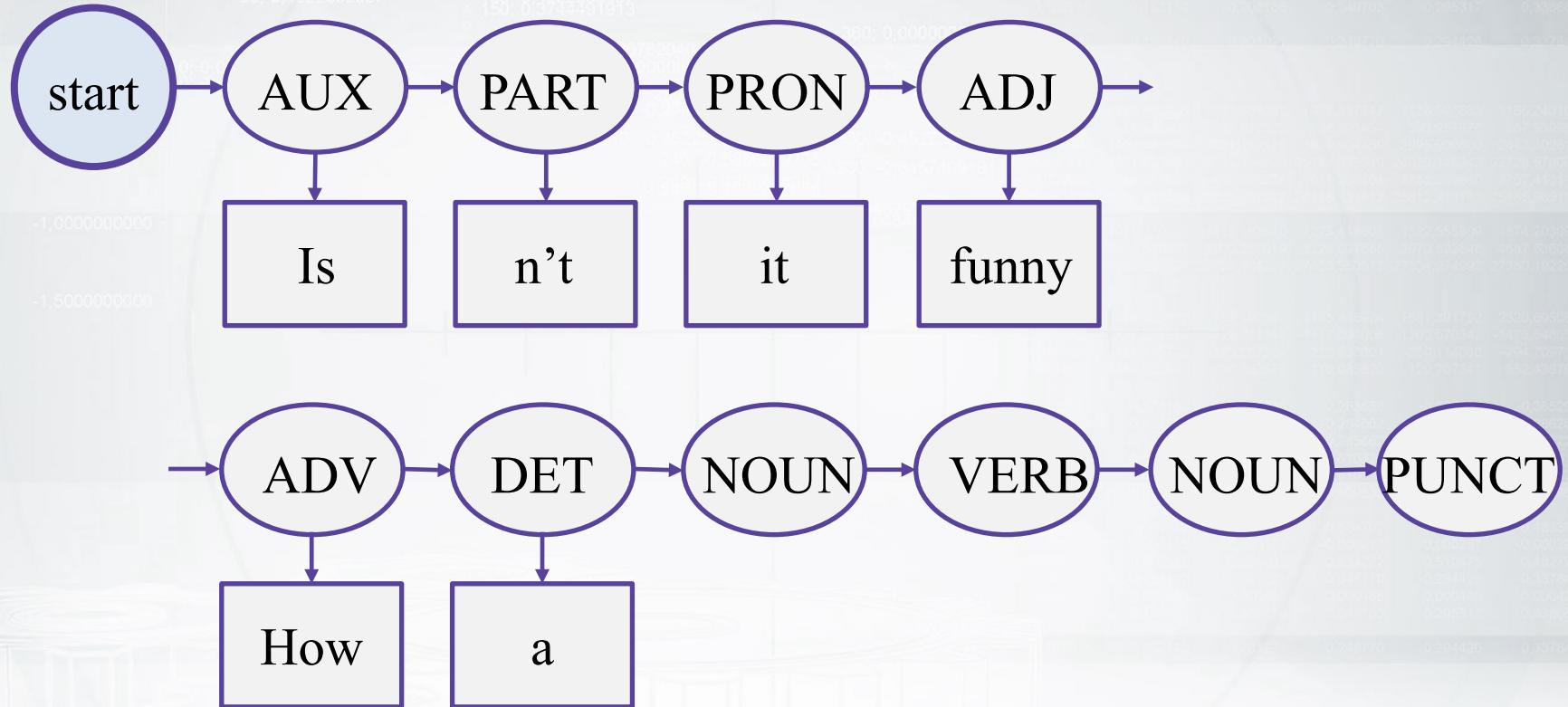
Text generation in HMM: an example



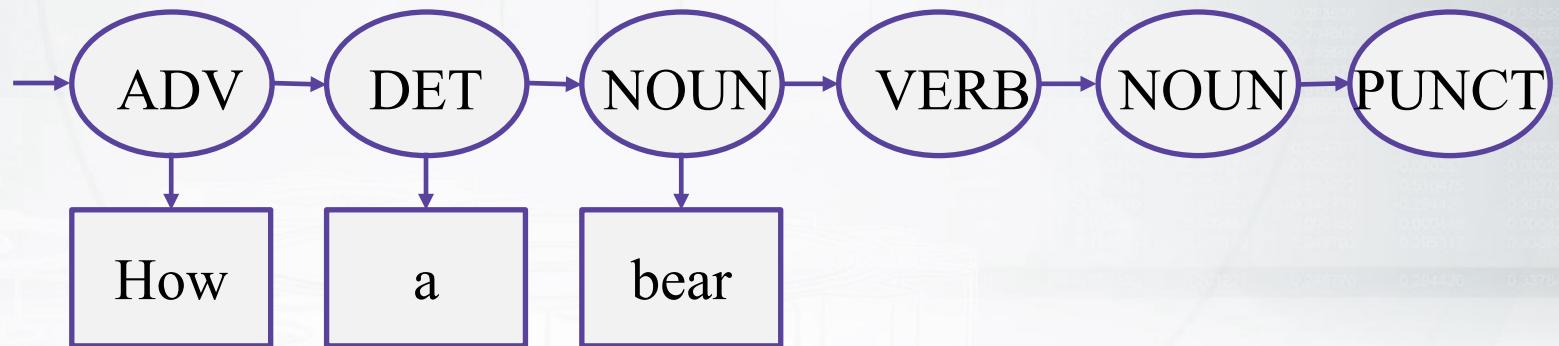
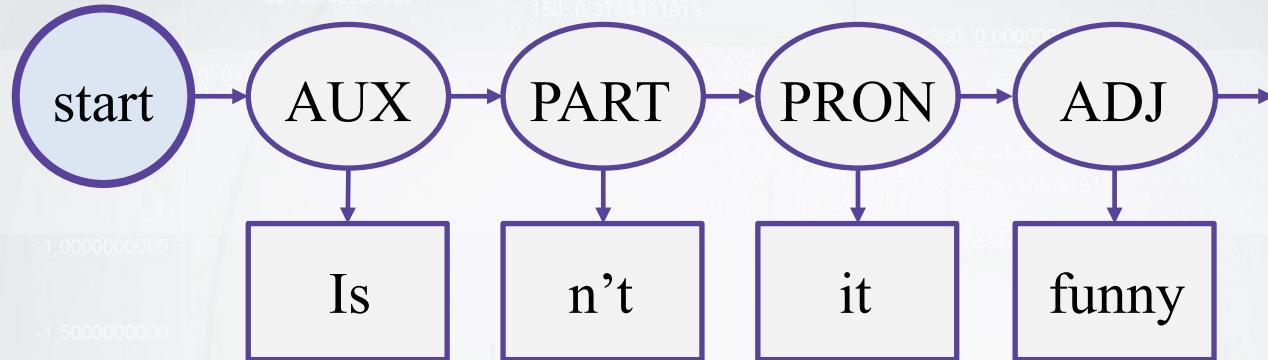
Text generation in HMM: an example



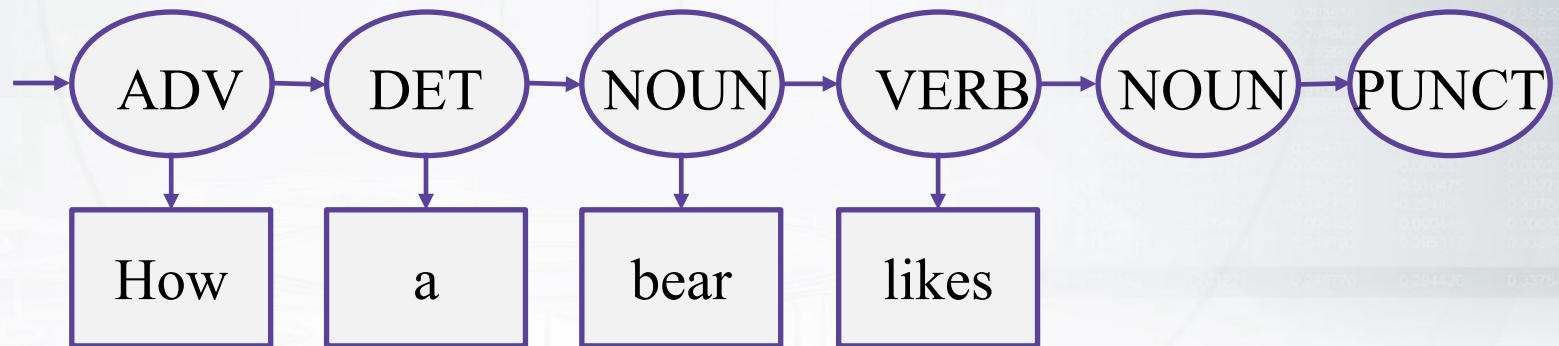
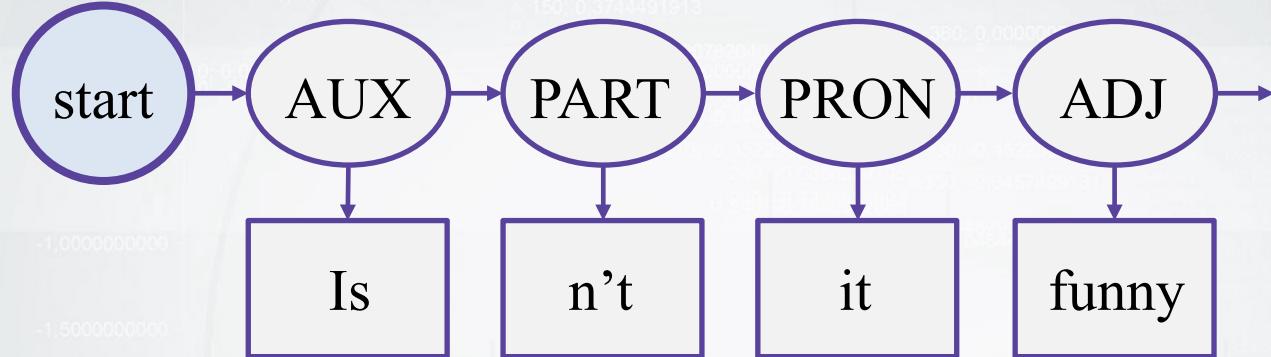
Text generation in HMM: an example



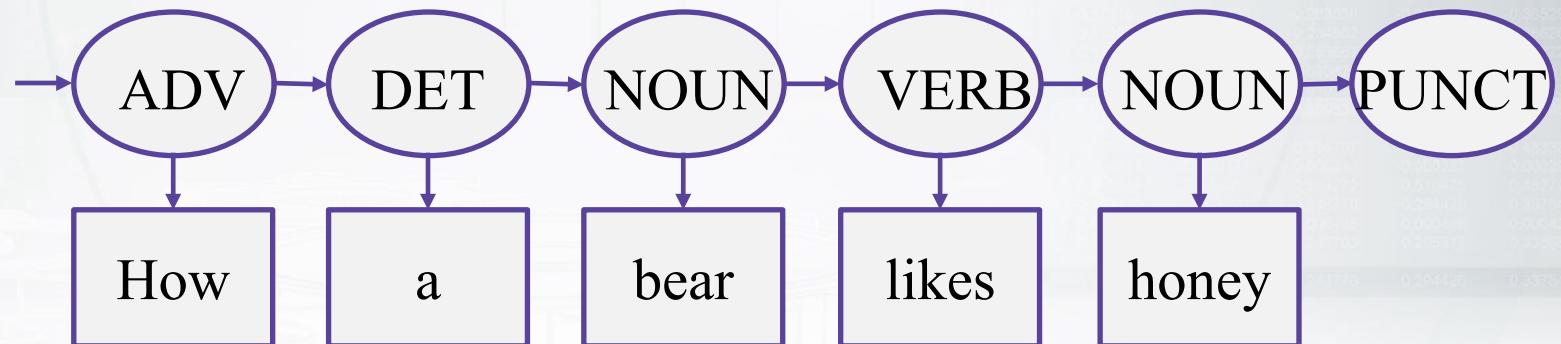
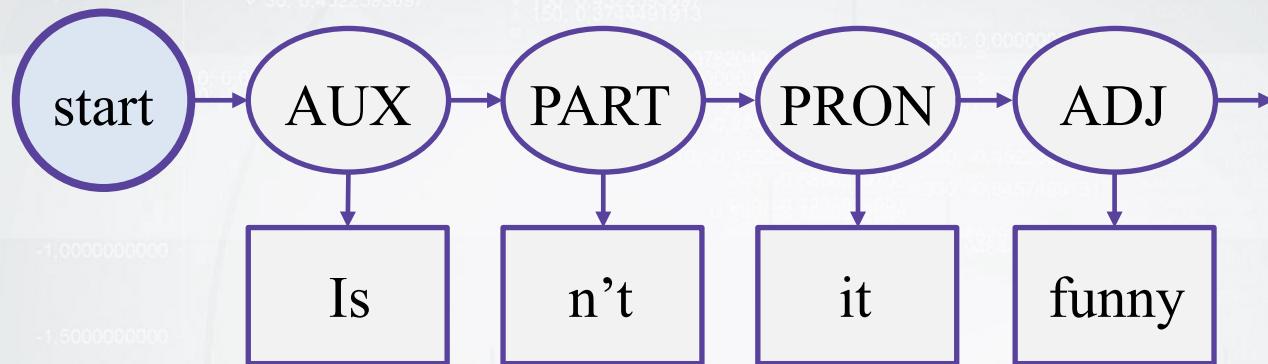
Text generation in HMM: an example



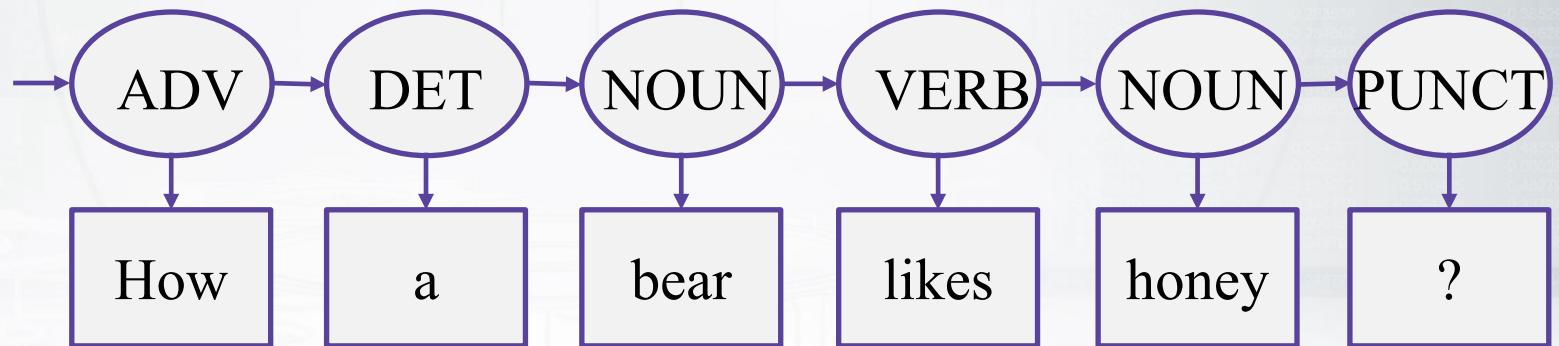
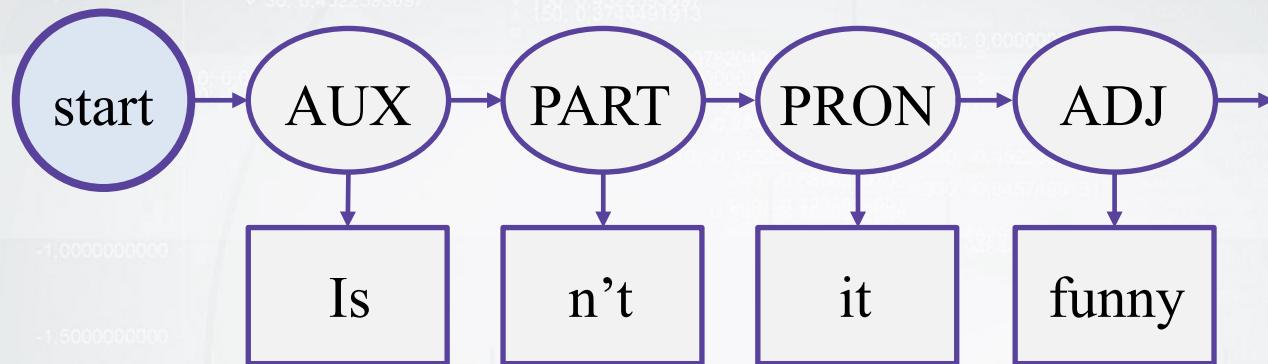
Text generation in HMM: an example



Text generation in HMM: an example



Text generation in HMM: an example

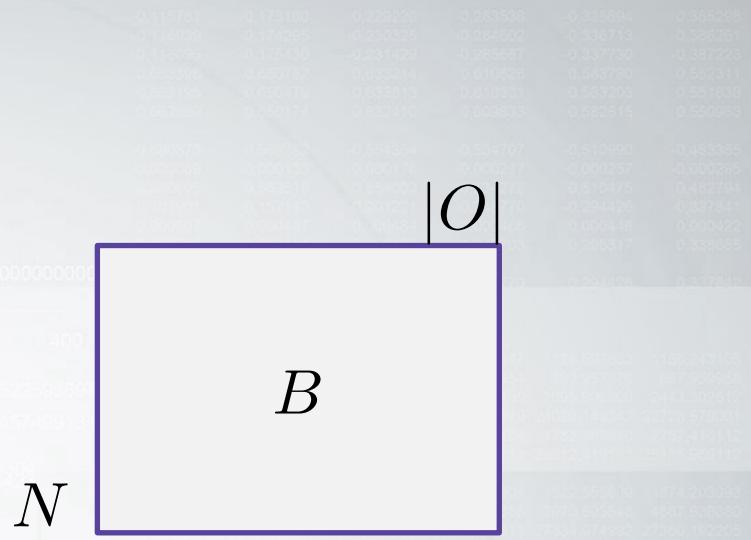
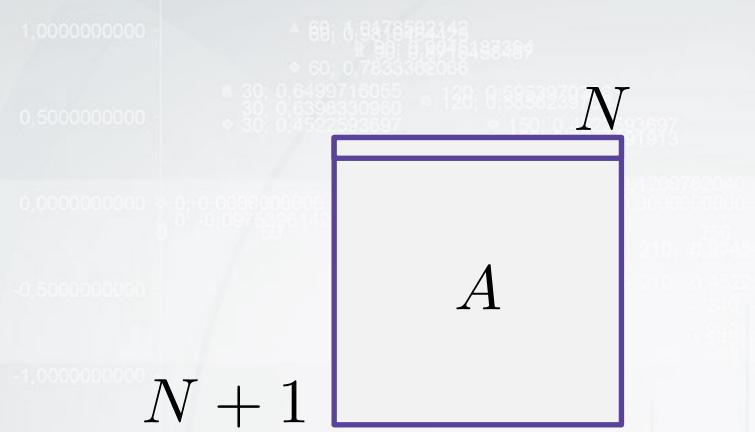


Formal definition of HMM

A Hidden Markov Model is specified by:

1. The set $S = s_1, s_2, \dots, s_N$ of hidden states
2. The start state s_0
3. The matrix A of transition probabilities: $a_{ij} = p(s_j | s_i)$
4. The set O of possible visible outcomes
5. The matrix B of output probabilities: $b_{kj} = p(o_k | s_j)$

How to train the model?



If we could see the tags in train set, we would count:

$$a_{ij} = p(s_j | s_i) = \frac{c(s_i \rightarrow s_j)}{c(s_i)}$$

$$b_{ik} = p(o_k | s_i) = \frac{c(s_i \rightarrow o_k)}{c(s_i)}$$

Supervised case: MLE

The same in more formal terms (this is MLE!):

$$a_{ij} = p(s_j | s_i) = \frac{\sum_{t=1}^T [y_{t-1} = s_i, y_t = s_j]}{\sum_{t=1}^T [y_t = s_i]}$$

Note that the corpus is considered as a single sequence of length T with special states between the sentences.

Supervised case: MLE

The same in more formal terms (this is MLE!):

$$a_{ij} = p(s_j | s_i) = \frac{\sum_{t=1}^T [y_{t-1} = s_i, y_t = s_j]}{\sum_{t=1}^T [y_t = s_i]}$$

Note that the corpus is considered as a single sequence of length T with special states between the sentences.

But we do not see the labels! How to train the model?

Baum-Welch algorithm (a sketch)

E-step: posterior probabilities for hidden variables:

$$p(y_{t-1} = s_i, y_t = s_j)$$

Can be effectively done with dynamic programming
(forward-backward algorithm)

M-step: maximum likelihood updates for the parameters:

$$a_{ij} = p(s_j | s_i) = \frac{\sum_{t=1}^T p(y_{t-1} = s_i, y_t = s_j)}{\sum_{t=1}^T p(y_t = s_i)}$$