

# CS918: LECTURE 8

Sequence Classification and Part-Of-Speech Tagging

---

Arkaitz Zubiaga, 29<sup>th</sup> October, 2018

## RECAP: WHAT IS TEXT CLASSIFICATION?

- Having as input:
  - A **text document**  $d$
  - A **set of categories**  $C=\{c_1, ..., c_m\}$
- The **text classification** task **outputs**:
  - Predicted class  $c^*$  that **document**  $d$  **belongs to**.

$$c^* \in C$$

## RECAP: TEXT CLASSIFICATION APPROACHES

- **Rule-based classifiers**, e.g. if email contains 'viagra' → spam
  - Significant **manual effort** involved.
- **Supervised classification**:
  - **Given**: a hand-labeled set of **document-class pairs**  
 $(d_1, c_1), (d_2, c_2), \dots, (d_m, c_m) \rightarrow$  classified into  $C = \{c_1, \dots, c_j\}$
  - The classifier **learns a model** that can **classify new documents into C**.

## RECAP: EVALUATION

- Binary classification ( $k = 2$ ):
  - Precision, recall and F1 score.
- Multiclass classification ( $k > 2$ ):
  - Accuracy, ratio of correct predictions, can be problematic when labels are skewed.
  - Better, precision, recall and F1 score per class, then:
    - Micro-averaged evaluation.
    - Macro-averaged evaluation.

## RECAP: ERROR ANALYSIS

- **Error analysis:** can help us find out where our classifier can do better.
- No magic formula for performing error analysis.
  - Look **where we are doing wrong**, what labels particularly.
  - Do our errors **have some common characteristics**? Can we infer a new feature from that?
  - Could our **classifier be favouring one of the classes** (e.g. the majority class)?

## LECTURE 8: CONTENTS

- Sequence Classification
- Sequence Classifiers:
  - Hidden Markov Models (HMM).
  - Maximum Entropy Markov Models (MEMM).
  - Conditional Random Fields (CRF).
- Using Sequence Classifiers for Part-of-Speech (POS) Tagging.

## SEQUENCE CLASSIFICATION

- Sometimes, **classification of items in a sequence is dependent on other sequence items:**
  - **Part-of-Speech (POS) tagging:**  
Assigning **categories to words**, e.g. adjective, noun or verb

Example:

The	man	saw	a	cat
<u>DET</u>	<u>NN</u>	<u>VB</u>	<u>DET</u>	<u>NN</u>

DET: determiner
NN: noun
VB: verb

## SEQUENCE CLASSIFICATION

- **Why** is classification **dependent** on other sequence items?

In our example, “The man saw a cat”:

**‘saw’ can be:**

a **noun** (if it refers to the tool for cutting)

a **verb** (if it’s simple past of ‘see’)

we can’t classify whether ‘saw’ is verb or noun by looking at the word alone → need to look at context, the sequence



## SEQUENCE CLASSIFICATION

<u>The</u>	<u>man</u>	<u>saw</u>	<u>a</u>	<u>cat</u>
???	???	???	???	???

- The first item, **'the'**, is **easy to classify**, it's always a determiner:

<u>The</u>	<u>man</u>	<u>saw</u>	<u>a</u>	<u>cat</u>
DET	???	???	???	???

## SEQUENCE CLASSIFICATION

The man saw a cat  
 DET    ???    ???    ???    ???

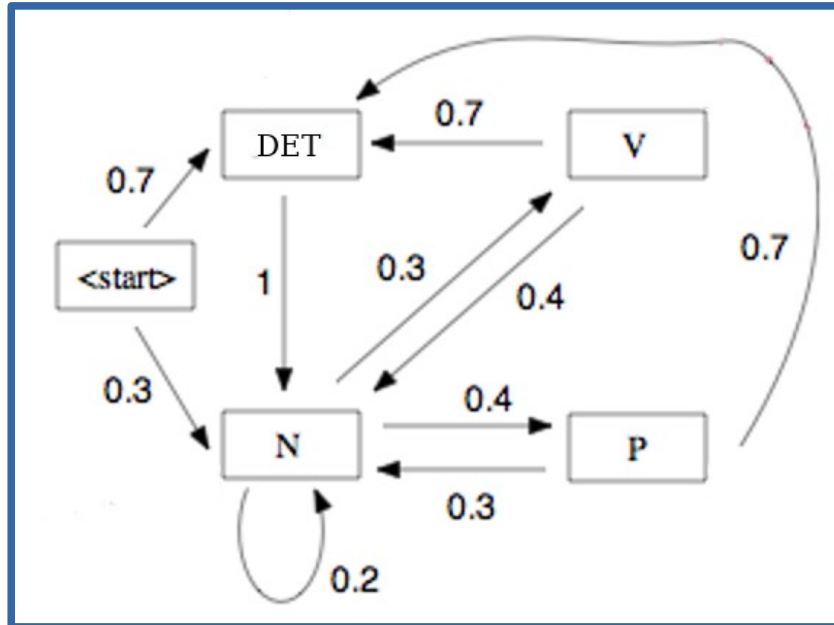
- Now **'man' is ambiguous**: (1) **noun**, i.e. male person, or (2) **verb**, i.e. 'take charge of'? We can look at:
  - $P(\text{'man}_{\text{NN}})$  vs  $P(\text{'man}_{\text{VB}})$   $\rightarrow$  probability of the word as noun or verb
  - $P(\text{NN} \mid \text{DET})$  vs  $P(\text{VB} \mid \text{DET})$   $\rightarrow$  probability of a noun or a verb to be preceded by a determiner

## SEQUENCE CLASSIFICATION

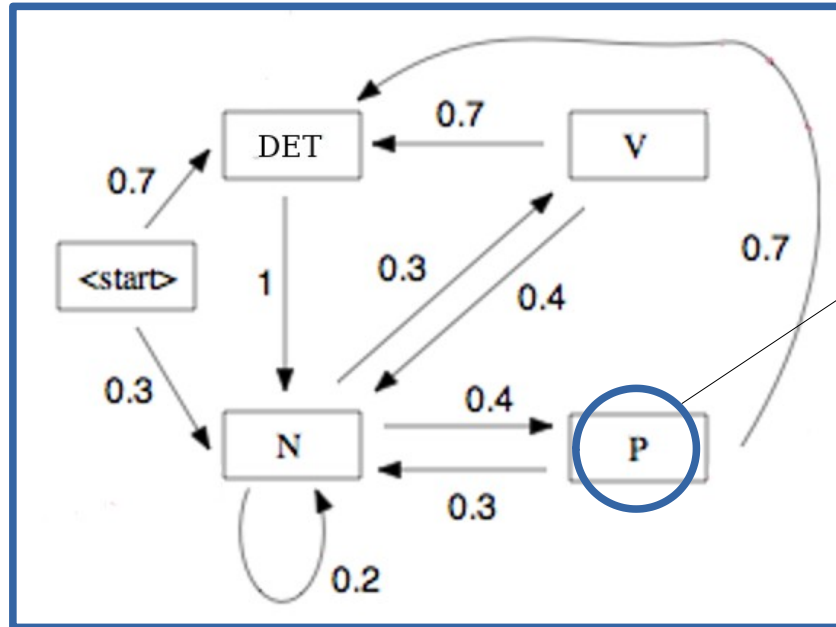
- **Two probabilities** to determine POS of a word:  
[1]  $P(\text{'manNN'})$  vs  $P(\text{'manVB'})$   
& [2]  $P(\text{'manNN'} \mid \text{DET})$  vs  $P(\text{'manVB'} \mid \text{DET})$
- Classifiers we have seen so far (NB, MaxEnt, SVM) can handle [1].
  - But they can't handle [2].

## SEQUENCE CLASSIFICATION

- Using training data, we can learn of word category  $POS_i$  being preceded by  $POS_j$ .

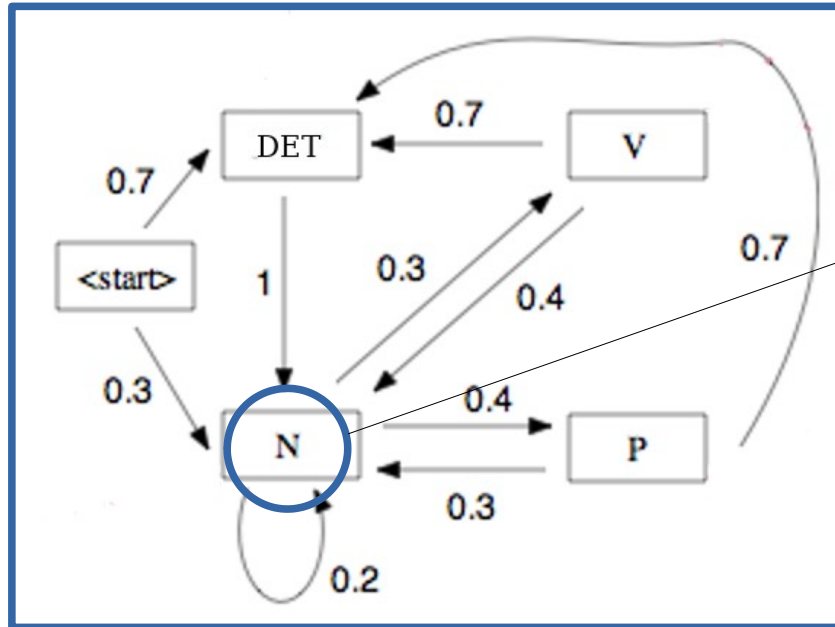


## SEQUENCE CLASSIFICATION



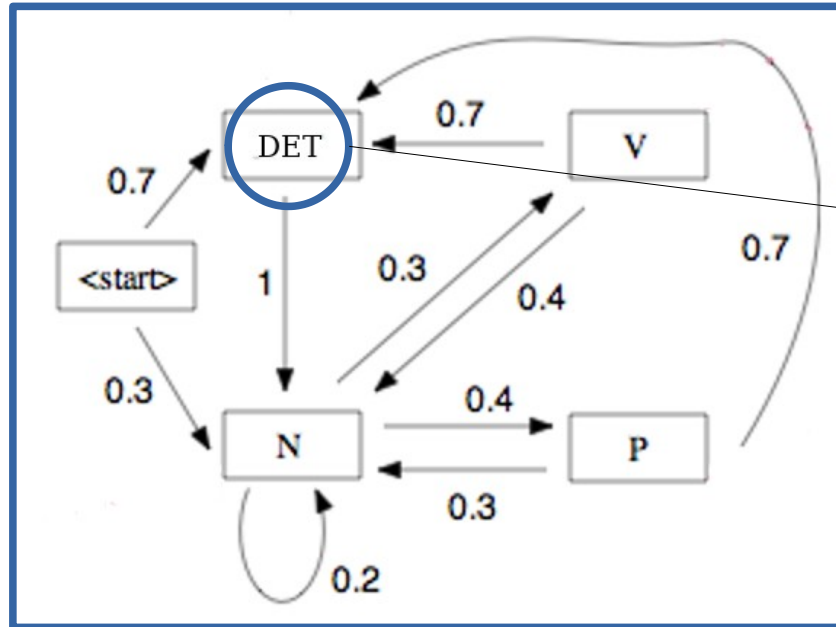
Prepositions can be followed by nouns (.7) or determiners (.3), never by verbs.

## SEQUENCE CLASSIFICATION



Nouns can be followed by verbs, prepositions, or another noun.

## SEQUENCE CLASSIFICATION



Determiners are ALWAYS followed by a noun.

## SEQUENCE CLASSIFICATION

[1]  $P(\text{'manNN'})$  vs  $P(\text{'manVB'})$

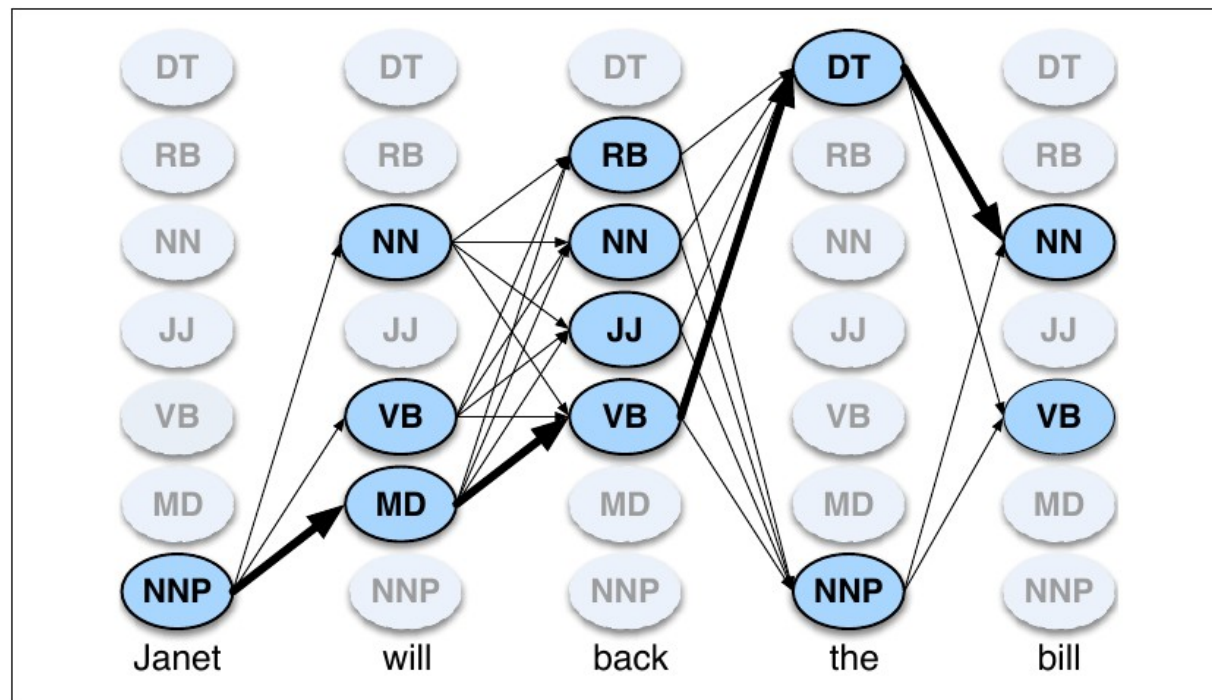
[2]  $P(\text{NN} \mid \text{DET})$  vs  $P(\text{VB} \mid \text{DET}) \rightarrow 0$

The	man	saw	a	cat
DET	???	???	???	???

No matter the outcome of [1], we know it's NN thanks to [2].



## FINDING THE MOST LIKELY SEQUENCE



## SEQUENCE CLASSIFICATION

- OK, and how do we achieve this?
  - Hidden Markov Models (HMM)
  - Maximum Entropy Markov Models (MEMM)
  - Conditional Random Fields (CRF)
  - Deep Learning:
    - Recurrent Neural Networks (RNN).
    - Long/Short-Term Memory Networks (LSTM).
    - Convolutional Neural Networks (CNN).

## SEQUENCE CLASSIFICATION

- OK, and how do we achieve this?
  - **Hidden Markov Models (HMM)**
  - **Maximum Entropy Markov Models (MEMM)**
  - **Conditional Random Fields (CRF)**
  - Deep Learning:
    - Recurrent Neural Networks (RNN).
    - Long/Short-Term Memory Networks (LSTM).
    - Convolutional Neural Networks (CNN).

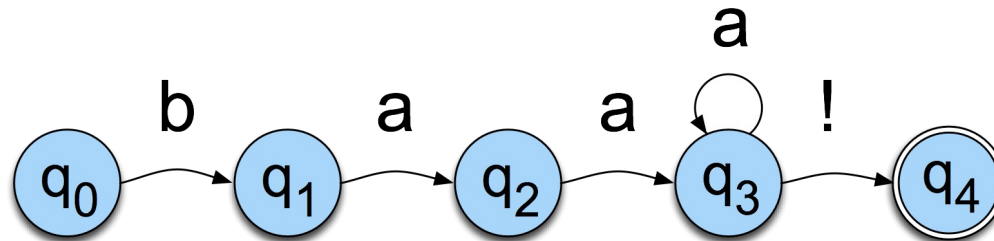


WARWICK

# HIDDEN MARKOV MODELS

## WHAT IS A MARKOV CHAIN?

- A Markov chain:
  - Set of **nodes connected with probabilities**, where weights on all edges leaving a node sum to 1.



# MARKOV CHAINS

$$Q = q_1 q_2 \dots q_N$$

a set of  $N$  **states**

$$A = a_{01} a_{02} \dots a_{n1} \dots a_{nm}$$

a **transition probability matrix**  $A$ , each  $a_{ij}$  representing the probability of moving from state  $i$  to state  $j$ , s.t.  $\sum_{j=1}^n a_{ij} = 1 \quad \forall i$

$$q_0, q_F$$

a special **start state** and **end (final) state** that are not associated with observations

## MARKOV CHAINS

- In a First order Markov Chain, the probability of a particular state depends ONLY on the previous state.

**Markov Assumption:**  $P(q_i | q_1 \dots q_{i-1}) \approx P(q_i | q_{i-1})$

- **Remember:** Lecture 3 on language models, Markov assumption:  
 $P(\text{library} \mid \text{I found two pounds in the}) \approx P(\text{library} \mid \text{the})$

## WHAT IS A HIDDEN MARKOV MODEL (HMM)?

$Q = q_1 q_2 \dots q_N$	a set of $N$ <b>states</b>
$A = a_{11} a_{12} \dots a_{n1} \dots a_{nn}$	a <b>transition probability matrix</b> $A$ , each $a_{ij}$ representing the probability of moving from state $i$ to state $j$ , s.t. $\sum_{j=1}^n a_{ij} = 1 \quad \forall i$
$O = o_1 o_2 \dots o_T$	a sequence of $T$ <b>observations</b> , each one drawn from a vocabulary $V = v_1, v_2, \dots, v_V$
$B = b_i(o_t)$	a sequence of <b>observation likelihoods</b> , also called <b>emission probabilities</b> , each expressing the probability of an observation $o_t$ being generated from a state $i$
$q_0, q_F$	a special <b>start state</b> and <b>end (final) state</b> that are not associated with observations, together with transition probabilities $a_{01} a_{02} \dots a_{0n}$ out of the start state and $a_{1F} a_{2F} \dots a_{nF}$ into the end state



## HMM: TWO KINDS OF PROBABILITIES

- **Transition probabilities:** pairwise probabilities of tags being preceded/followed by another tag.  
e.g. determiner likely followed by noun or adjective.
- **Emission probabilities:** probability for a particular tag to be a particular word, e.g. verb is very likely to be “is”.

## TRANSITION PROBABILITIES: EXAMPLE

	<b>NNP</b>	<b>MD</b>	<b>VB</b>	<b>JJ</b>	<b>NN</b>	<b>RB</b>	<b>DT</b>
<b>&lt;s&gt;</b>	0.2767	0.0006	0.0031	0.0453	0.0449	0.0510	0.2026
<b>NNP</b>	0.3777	0.0110	0.0009	0.0084	0.0584	0.0090	0.0025
<b>MD</b>	0.0008	0.0002	0.7968	0.0005	0.0008	0.1698	0.0041
<b>VB</b>	0.0322	0.0005	0.0050	0.0837	0.0615	0.0514	0.2231
<b>JJ</b>	0.0366	0.0004	0.0001	0.0733	0.4509	0.0036	0.0036
<b>NN</b>	0.0096	0.0176	0.0014	0.0086	0.1216	0.0177	0.0068
<b>RB</b>	0.0068	0.0102	0.1011	0.1012	0.0120	0.0728	0.0479
<b>DT</b>	0.1147	0.0021	0.0002	0.2157	0.4744	0.0102	0.0017



## EMISSION PROBABILITIES: EXAMPLE

	<b>Janet</b>	<b>will</b>	<b>back</b>	<b>the</b>	<b>bill</b>
<b>NNP</b>	0.000032	0	0	0.000048	0
<b>MD</b>	0	0.308431	0	0	0
<b>VB</b>	0	0.000028	0.000672	0	0.000028
<b>JJ</b>	0	0	0.000340	0	0
<b>NN</b>	0	0.000200	0.000223	0	0.002337
<b>RB</b>	0	0	0.010446	0	0
<b>DT</b>	0	0	0	0.506099	0

## TWO IMPORTANT ASSUMPTIONS IN HMM

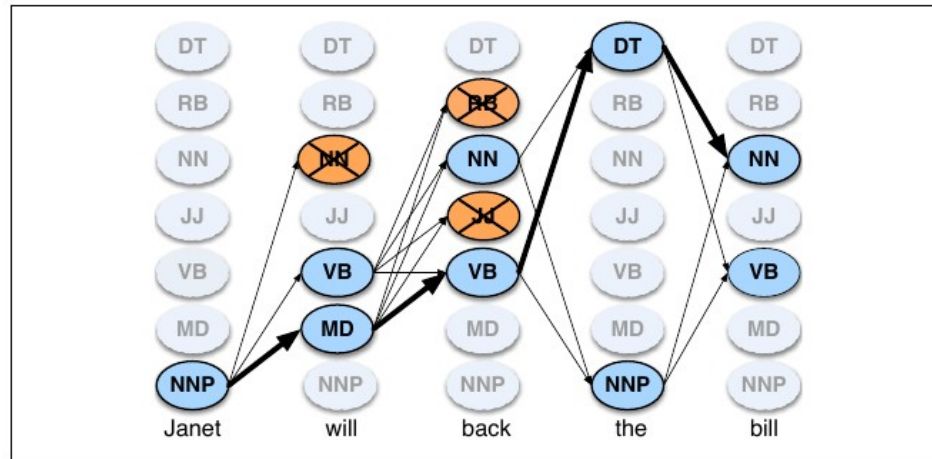
- Again, the Markov assumption, i.e. dependence only on the previous state.

$$\text{Markov Assumption: } P(q_i | q_1 \dots q_{i-1}) \approx P(q_i | q_{i-1})$$

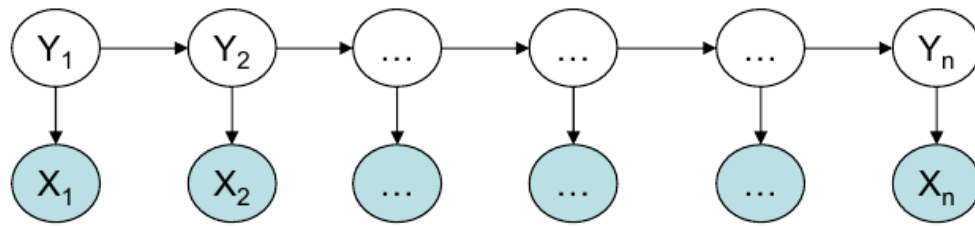
- There must be a dependency between sequence items, e.g.:
  -  Weather (dependency): this morning's weather may determine the probability of this afternoon's weather.
  -  Independent events: my laptop having been broken doesn't determine the probability of my next laptop breaking.

## OPTIMISING WITH BEAM SEARCH

- Beam search: only consider top  $\beta$  options for each word.
- Here  $\beta = 2$ :



## LIMITATIONS OF HMM

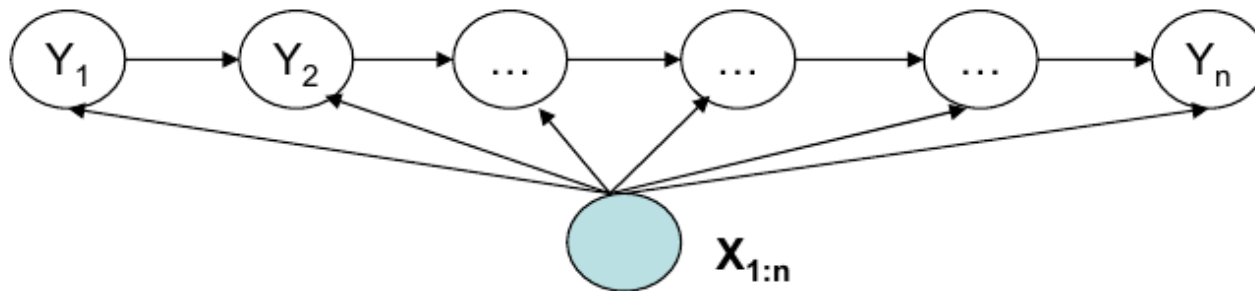


- Limitations of HMM:
  - Models dependencies between each state and only its corresponding observation.
  - Learns joint distribution of states and observations  $P(Y, X)$ , but not the conditional probability  $P(Y|X)$ .

## HMM IMPLEMENTATIONS

- NLTK HMM:  
[http://www.nltk.org/\\_modules/nltk/tag/hmm.html](http://www.nltk.org/_modules/nltk/tag/hmm.html)
- hmmlearn:  
<https://github.com/hmmlearn/hmmlearn>
- Scikit HMM:  
<http://scikit-learn.sourceforge.net/stable/modules/hmm.html>

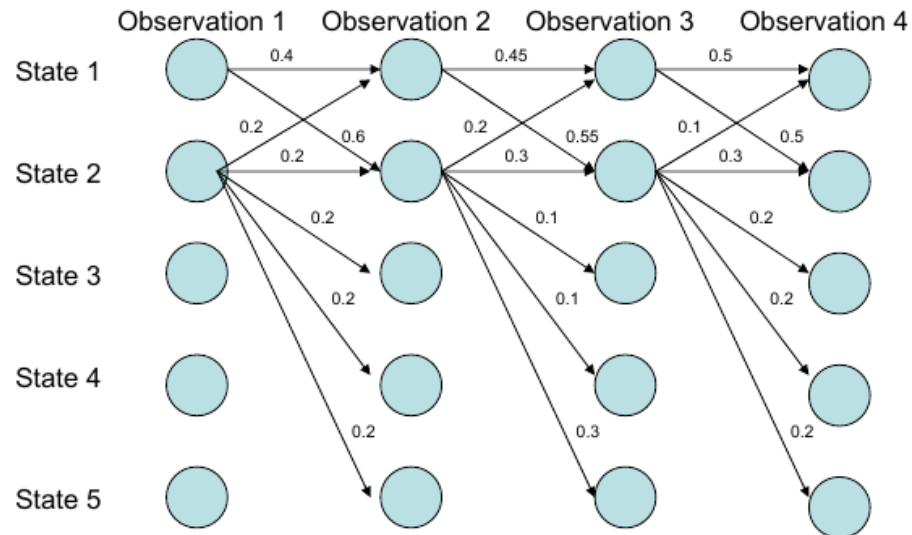
## MEMM: ALTERNATIVE TO HMM



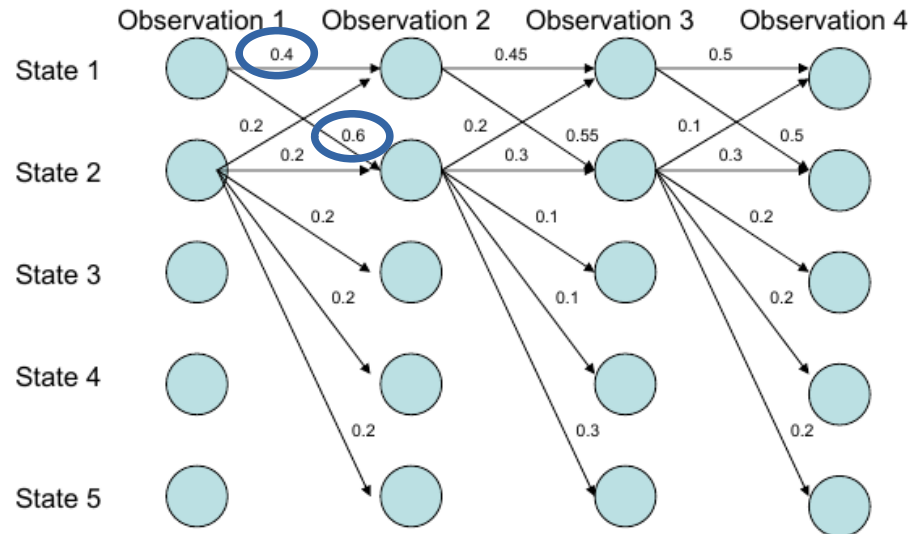
- Maximum Entropy Markov Models (MEMM): Models **dependencies** between **each state and the full observation sequence**.
- Learning objective function consistent with predictive function:  $P(Y|X)$ .



## MEMM: THE LABEL BIAS PROBLEM

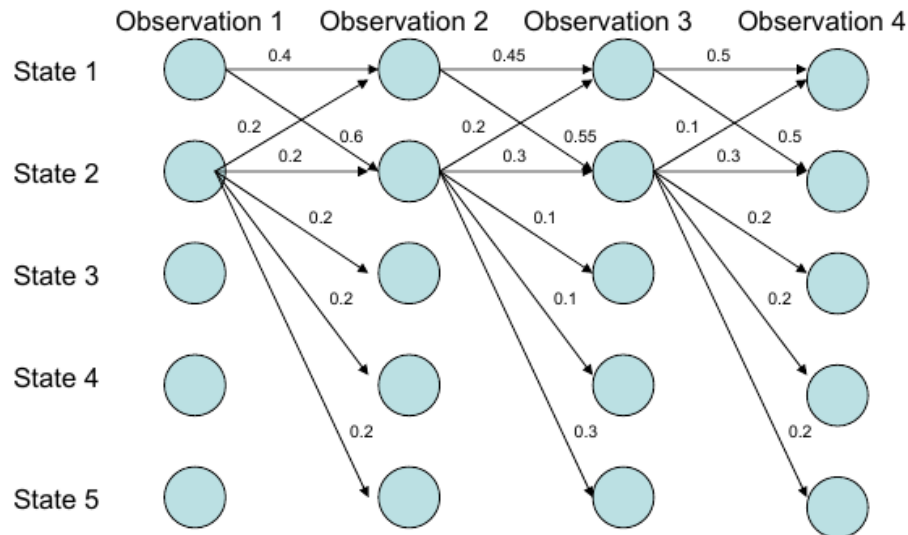


## MEMM: THE LABEL BIAS PROBLEM



- In principle, locally, state 1 prefers state 2 in 2<sup>nd</sup> place.

## MEMM: THE LABEL BIAS PROBLEM



Despite  $1 \rightarrow 2$  being more likely **locally**,

the **entire path probability** will favour  $1 \rightarrow 1$

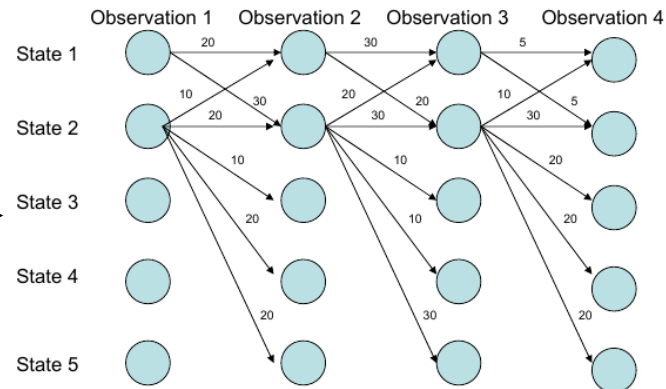
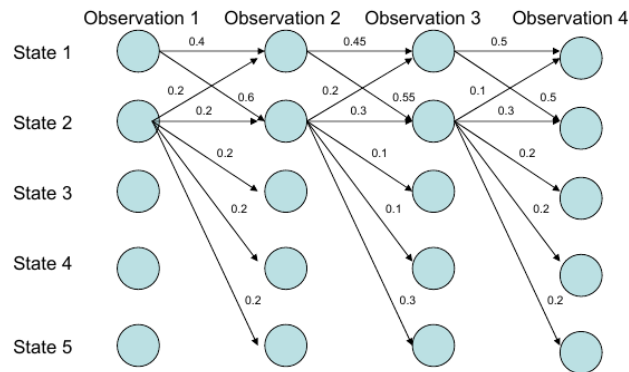
- But if we look at the entire path (which MEMM does):

- $P(1 \rightarrow 1 \rightarrow 1 \rightarrow 1) = 0.4 * 0.45 * 0.5 = 0.09$

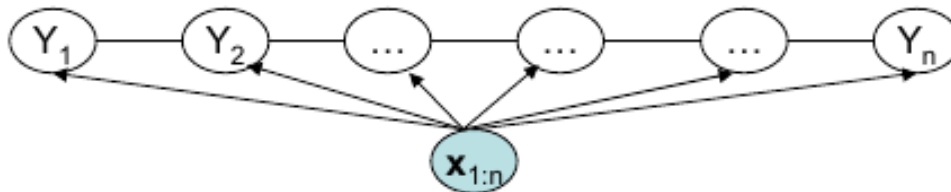
- $P(1 \rightarrow 2 \rightarrow 2 \rightarrow 2) = 0.6 * 0.3 * 0.3 = 0.054$

## SOLUTION TO LABEL BIAS

- Solution: don't normalise probabilities locally, use local potentials.



## CONDITIONAL RANDOM FIELDS



$$P(\mathbf{y}_{1:n} | \mathbf{x}_{1:n}) = \frac{1}{Z(\mathbf{x}_{1:n})} \prod_{i=1}^n \phi(y_i, y_{i-1}, \mathbf{x}_{1:n}) = \frac{1}{Z(\mathbf{x}_{1:n}, \mathbf{w})} \prod_{i=1}^n \exp(\mathbf{w}^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_{1:n}))$$

- CRF:
  - Global normaliser  $Z(\mathbf{x})$  overcomes label bias issue of MEMM.
  - Models the dependency between each state and the entire observation sequence (like MEMM).

## CONDITIONAL RANDOM FIELDS

- Widely used for a range of sequence classification tasks:

Lafferty, J., McCallum, A., & Pereira, F. C. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data.

[https://repository.upenn.edu/cgi/viewcontent.cgi?article=1162&context=cis\\_papers](https://repository.upenn.edu/cgi/viewcontent.cgi?article=1162&context=cis_papers)

## IMPLEMENTATIONS OF CRF

- Python-crfsuite:  
<http://python-crfsuite.readthedocs.io/en/latest/>
- PyStruct:  
<https://pystruct.github.io/>

## BEYOND CLASSIFIER LINEAR SEQUENCES

- We've talked about classifier linear sequences.
- But sometimes we can have more complex sequences:
  - Graph or tree-structured sequences.



## WHERE DO WE FIND TREE SEQUENCES IN NLP?

- Post classification in forums, e.g. binary classification: does a post answer the question of the 1<sup>st</sup> post?
- Post 1
  - Post 1.1
    - Post 1.1.1
    - Post 1.1.2
  - Post 1.2
    - Post 1.2.1
      - Post 1.2.1.1
    - Post 1.2.2

# HMM VS CRF

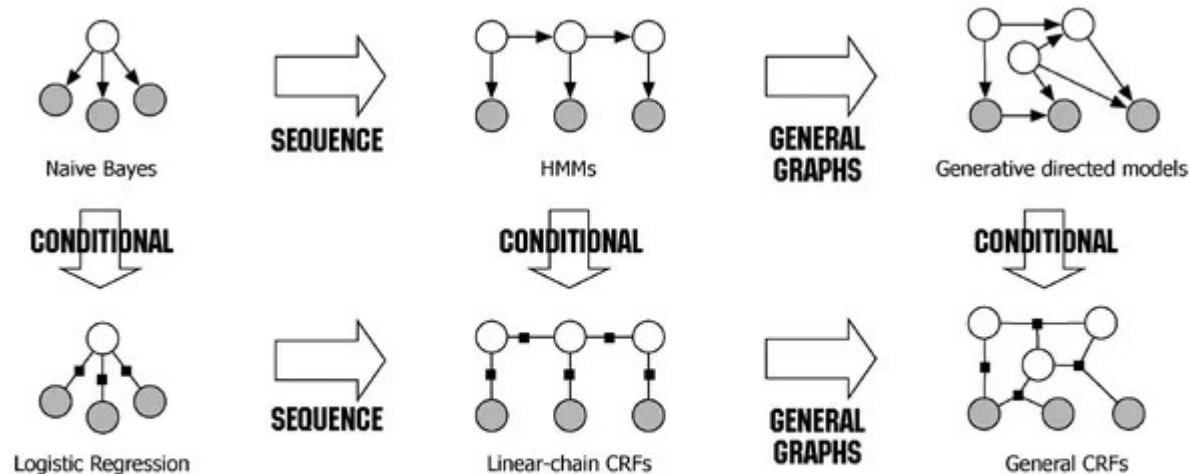


Fig. 2.4 Diagram of the relationship between naive Bayes, logistic regression, HMMs, linear-chain CRFs, generative models, and general CRFs.



# USING SEQUENCE CLASSIFIERS FOR PART-OF-SPEECH (POS) TAGGING

## PART-OF-SPEECH TAGGING

- As we were saying, sequence can play an important role in determining POS tags in a sentence:

The	man	saw	a	cat
DET	NN	VB	DET	NN

“man” can’t be a verb if it’s preceded by a determiner.

## PART-OF-SPEECH TAGGING

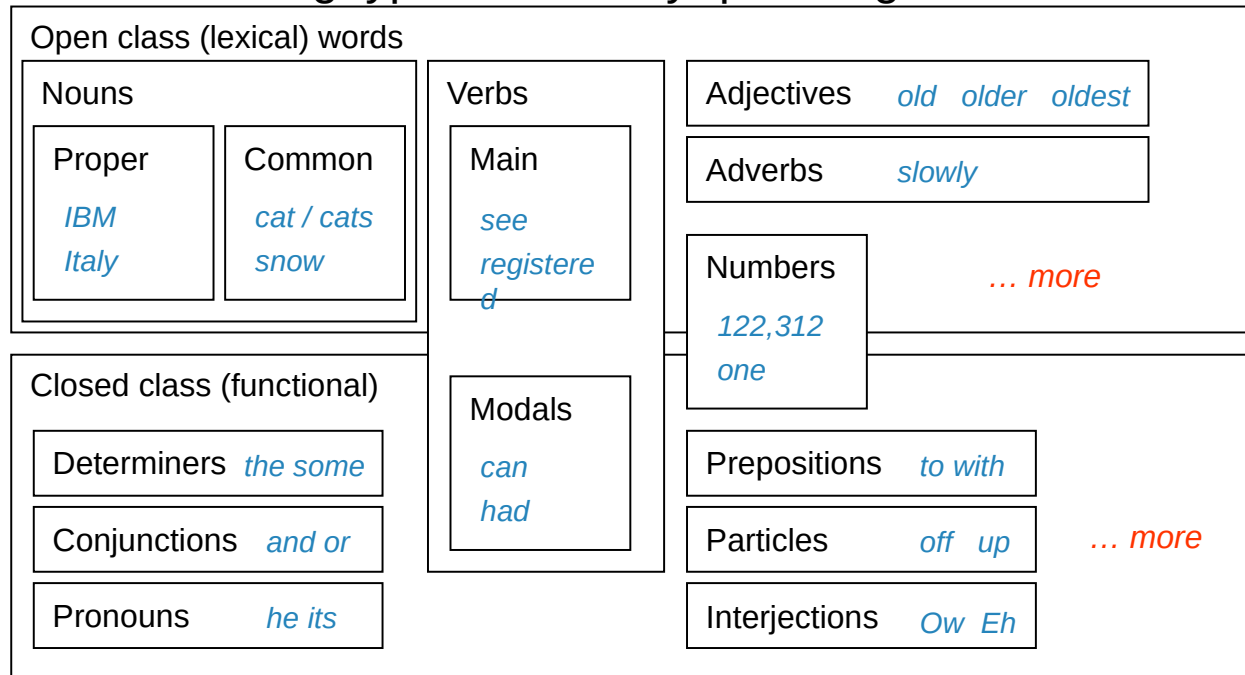
- As we were saying, sequence can play an important role in determining POS tags in a sentence:

The	man	saw	a	cat
DET	NN	VB	DET	NN

- With **HMM**, only looking at previous label, we **can end up predicting sequence of 5 nouns** (NN, NN, NN, NN, NN).
- Looking at the **entire sequence** (MEMM, CRF), we avoid this, **we never have a sentence only made of 5 nouns**.

## PART-OF-SPEECH TAGGING

- The list of POS tag types is actually quite large!



# THE PENN TREEBANK TAG SET

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coordinating conjunction	<i>and, but, or</i>	PDT	predeterminer	<i>all, both</i>	VBP	verb non-3sg present	<i>eat</i>
CD	cardinal number	<i>one, two</i>	POS	possessive ending	<i>'s</i>	VBZ	verb 3sg pres	<i>eats</i>
DT	determiner	<i>a, the</i>	PRP	personal pronoun	<i>I, you, he</i>	WDT	wh-determ.	<i>which, that</i>
EX	existential 'there'	<i>there</i>	PRP\$	possess. pronoun	<i>your, one's</i>	WP	wh-pronoun	<i>what, who</i>
FW	foreign word	<i>mea culpa</i>	RB	adverb	<i>quickly</i>	WP\$	wh-possess.	<i>whose</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	RBR	comparative adverb	<i>faster</i>	WRB	wh-adverb	<i>how, where</i>
JJ	adjective	<i>yellow</i>	RBS	superlatv. adverb	<i>fastest</i>	\$	dollar sign	<i>\$</i>
JJR	comparative adj	<i>bigger</i>	RP	particle	<i>up, off</i>	#	pound sign	<i>#</i>
JJS	superlative adj	<i>wildest</i>	SYM	symbol	<i>+, %, &amp;</i>	“	left quote	<i>' or “</i>
LS	list item marker	<i>1, 2, One</i>	TO	“to”	<i>to</i>	”	right quote	<i>' or ”</i>
MD	modal	<i>can, should</i>	UH	interjection	<i>ah, oops</i>	(	left paren	<i>[, (, {, &lt;</i>
NN	sing or mass noun	<i>llama</i>	VB	verb base form	<i>eat</i>	)	right paren	<i>], ), }, &gt;</i>
NNS	noun, plural	<i>llamas</i>	VBD	verb past tense	<i>ate</i>	,	comma	<i>,</i>
NNP	proper noun, sing.	<i>IBM</i>	VBG	verb gerund	<i>eating</i>	.	sent-end punc	<i>. ! ?</i>
NNPS	proper noun, plu.	<i>Carolinas</i>	VBN	verb past part.	<i>eaten</i>	:	sent-mid punc	<i>: ; ... - -</i>

## OPEN VS CLOSED CLASSES

- Open vs. Closed classes:
  - **Closed:**
    - **determiners:** a, an, the,...
    - **pronouns:** she, he, I,...
    - **prepositions:** on, under, over, near, by,...
  - **Open:**
    - **Nouns, Verbs, Adjectives, Adverbs.**



## CHALLENGES IN POS TAGGING: AMBIGUITY

- Words often have more than one possible POS, e.g. **back**:
  - The **back** door = JJ (adjective)
  - On my **back** = NN (noun)
  - Win the voters **back** = RB (adverb)
  - Promised to **back** the bill = VB (verb)
- See list of POS tags:  
<http://rednoise.org/rita/reference/PennTags.html>

## PART-OF-SPEECH TAGGING

- POS tagging example:
  - Input: **Plays well with others**

NNS UH IN NNS  
VBZ JJ  
NN  
RB

1. List candidate labels for each word.
2. Based on probabilities learnt from training data, the classifier predicts the most likely POS sequence.

NB: the fact that there are 2 unambiguous words (with & others) is useful to first label them, and then predict the other 2.

- Output: **Plays/VBZ well/RB with/IN others/NNS**

## PART-OF-SPEECH TAGGING

- **Uses of POS tagging** in NLP:
  - Text-to-speech.
  - Phrase search through regular expressions, e.g. (Det) Adj\* N+
  - As input to or to speed up a full linguistic parser (later lectures)
  - If you know the tag, you can back off to it, e.g. in lemmatisation, saw → see, or saw → saw?
- In other fields:
  - Computational biology.
  - Prediction of series of data, e.g. weather forecasting.

## POS TAGGING PERFORMANCE

- Current POS tagging system are very accurate:
  - **Baseline** system that predicts the most common POS for a word already gets **90% accuracy** → thanks to many words not being ambiguous.
  - **Standard classifiers** (no sequence) can achieve **93% accuracy**.
  - **Sequence classifiers** can achieve **97% accuracy**.

## POS TAGGING: AMBIGUITY

<b>Types:</b>		<b>WSJ</b>	<b>Brown</b>
<b>Unambiguous</b>	(1 tag)	44,432 ( <b>86%</b> )	45,799 ( <b>85%</b> )
<b>Ambiguous</b>	(2+ tags)	7,025 ( <b>14%</b> )	8,050 ( <b>15%</b> )
<b>Tokens:</b>			
<b>Unambiguous</b>	(1 tag)	577,421 ( <b>45%</b> )	384,349 ( <b>33%</b> )
<b>Ambiguous</b>	(2+ tags)	711,780 ( <b>55%</b> )	786,646 ( <b>67%</b> )

## REFERENCES

- List of software for POS tagging:  
[https://aclweb.org/aclwiki/Part-of-speech\\_tagging](https://aclweb.org/aclwiki/Part-of-speech_tagging)

## ASSOCIATED READING

- Jurafsky, Daniel, and James H. Martin. 2009. Speech and Language Processing: An Introduction to Natural Language Processing, Speech Recognition, and Computational Linguistics. 3rd edition. **Chapters 8.**
- Bird Steven, Ewan Klein, and Edward Loper. Natural Language Processing with Python. O'Reilly Media, Inc., 2009. **Chapter 6 Section 1.6.**