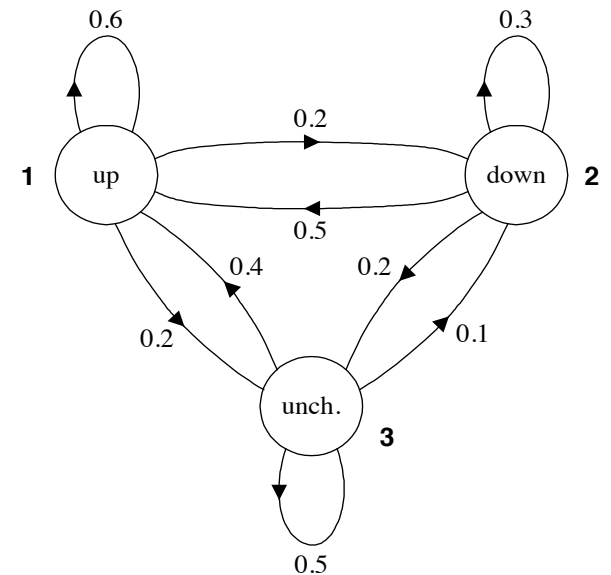


# Sequence Tagging: hidden Markov models

COMP90042 Lecture 15



THE UNIVERSITY OF  
MELBOURNE



# Pos tagging recap

- Janet will back the bill
- Janet/**NNP** will/**MB** back/**VP** the/**DT** bill/**NN**
- Local classifier: prone to **error propagation**
- What about treating the full sequence as a “class”?
  - \* Output: “NNP\_MB\_VP\_DT\_NN”
- Problems:
  - \* Exponentially many combinations:  $|\text{Tags}|^M$ , where M is the length
  - \* How to tag sequences of different lengths?

# A better approach

- Tagging is a sentence-level task but as humans we **decompose** it into small word-level tasks.
  - \* Janet/**NNP** will/**MB** back/**VP** the/**DT** bill/**NN**
- Solution:
  - \* Define a model that decompose a tagging sentence into individual words
  - \* But that takes into account the whole sequence when learning and predicting (no error propagation)
- This is the idea of **sequence labelling**, and more general, **structured prediction**.

# A probabilistic model

- Goal: obtain best tag sequence  $\mathbf{t}$  from sentence  $\mathbf{w}$ 
  - \*  $\hat{\mathbf{t}} = \operatorname{argmax}_{\mathbf{t}} P(\mathbf{t}|\mathbf{w})$
  - \*  $\hat{\mathbf{t}} = \operatorname{argmax}_{\mathbf{t}} \frac{P(\mathbf{w}|\mathbf{t})P(\mathbf{t})}{P(\mathbf{w})} = \operatorname{argmax}_{\mathbf{t}} P(\mathbf{w}|\mathbf{t}) P(\mathbf{t})$   
[Bayes]
- Let's decompose:
  - \*  $P(\mathbf{w}|\mathbf{t}) = \prod_{i=1}^n P(w_i|t_i)$  [Prob. of a word depends only on the tag]
  - \*  $P(\mathbf{t}) = \prod_{i=1}^n P(t_i|t_{i-1})$  [Prob. of a tag depends only on the previous tag]
- These are **independence assumptions** (remember Naïve Bayes?)
- This is a Hidden Markov Model (HMM)

# Hidden Markov model

- $\hat{\mathbf{t}} = \operatorname{argmax}_{\mathbf{t}} P(\mathbf{w}|\mathbf{t}) P(\mathbf{t})$
- $P(\mathbf{w}|\mathbf{t}) = \prod_{i=1}^n P(w_i|t_i)$
- $P(\mathbf{t}) = \prod_{i=1}^n P(t_i|t_{i-1})$
- Why “Markov”?
  - \* Because it assumes the sequence follows a Markov chain: probability of an event (tag) depends only on the previous one (previous tag)
- Why “Hidden”?
  - \* Because the events (tags) are not seen: goal is to find the best sequence

# HMMs - training

- Parameters are the individual probabilities  $P(w_i|t_i)$  and  $P(t_i|t_{i-1})$ 
  - \* Respectively, **emission** and **transition** probabilities
- Training uses Maximum Likelihood Estimation (MLE)
  - \* In Naïve Bayes, this is done by simply counting word frequencies according to the class.
- We do **exactly the same** in HMMs!
  - \*  $P(\textit{like}|\textit{VB}) = \frac{\textit{count}(\textit{VB},\textit{like})}{\textit{count}(\textit{VB})}$
  - \*  $P(\textit{NN}|\textit{DT}) = \frac{\textit{count}(\textit{DT},\textit{NN})}{\textit{count}(\textit{DT})}$

# HMMs - training

- What about the first tag?
- Assume we have a symbol “<s>” that represents the start of your sentence.
  - \*  $P(NN | <s>) = \frac{\text{count}(<s>, NN)}{\text{count}(<s>)}$
- What about unseen (word,tag) and (tag, previous) combinations?
- Same as Naïve Bayes
  - \* Smoothing techniques

# Transition Matrix

	<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

**Figure 10.5** The  $A$  transition probabilities  $P(t_i|t_{i-1})$  computed from the WSJ corpus without smoothing. Rows are labeled with the conditioning event; thus  $P(VB|MD)$  is 0.7968.



# Emission (observation) Matrix

	<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	0.000097	0
<b>NN</b>	0	0.000200	0.000223	0.000006	0.002337
<b>RB</b>	0	0	0.010446	0	0
<b>DT</b>	0	0	0	0.506099	0

**Figure 10.6** Observation likelihoods  $B$  computed from the WSJ corpus without smoothing.

# HMMs – prediction (decoding)

$$\hat{t} = \operatorname{argmax}_t P(\mathbf{w}|\mathbf{t}) P(\mathbf{t})$$

$$= \operatorname{argmax}_t \prod_{i=1}^n P(w_i|t_i)P(t_i|t_{i-1})$$

- Simple idea: for each word, take the tag that maximises  $P(w_i|t_i)P(t_i|t_{i-1})$ . Do it left-to-right.
- This is wrong! We are looking for  $\operatorname{argmax}_t$ , not individual  $\operatorname{argmax}_{t_i}$  terms.
  - \* This is a local classifier: error propagation
- Correct way: take **all** possible tag combinations, evaluate them, take the max (like Naïve Bayes)
  - \* Problem: exponential number of sequences.

# The Viterbi algorithm

- Dynamic Programming to the rescue!
  - \* We can still proceed sequentially, as long as we careful.
- “a cat” -> a/**DT** cat/**NN**
- Best tag for “a” is easy: take  $\operatorname{argmax}_t P(a|t)P(t| < s >)$ 
  - \* We can do that because first “tag” is always “<s>”
- Suppose best tag for “a” is DT. To get the tag for “cat”, we can take  $\operatorname{argmax}_t P(cat|t)P(t|DT)$  but this is wrong.
- Instead, we keep track of **scores for each tag** for “a” and check **what would happen** if “a” had a different tag.

# The viterbi algorithm

	Janet	will	back	the	bill
NNP					
MD					
VB					
JJ					
NN					
RB					
DT					


# The viterbi algorithm

	<b>Janet</b>	<b>will</b>	<b>back</b>	<b>the</b>	<b>bill</b>
NNP	$P(\text{Janet}   \text{NNP}) * P(\text{NNP}   <s>)$				
MD	$P(\text{Janet}   \text{MD}) * P(\text{MD}   <s>)$				
VB	...				
JJ	...				
NN	...				
RB	...				
DT	...				

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	0.000032 * 0.2767				
MD	0 * 0.0006				
VB	...				
JJ	...				
NN	...				
RB	...				
DT	...				

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 				
MD	0				
VB	0				
JJ	0				
NN	0				
RB	0				
DT	0				

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 	$P(\text{will}   \text{NNP}) * P(\text{NNP}   t_{\text{Janet}}) * s(t_{\text{Janet}}   \text{Janet})$			
MD	0	...			
VB	0	...			
JJ	0	...			
NN	0	...			
RB	0	...			
DT	0	...			



# The viterbi algorithm


	Janet	will	back	the	bill
NNP	8.8544e-06 ●	$P(\text{will}   \text{NNP}) * P(\text{NNP}   t_{\text{Janet}}) * S(t_{\text{Janet}}   \text{Janet})$			
MD	0	...			
VB	0	...			
JJ	0	...			
NN	0	...			
RB	0	...			
DT	0	...			

Calculate this for all tags,  
take the max.

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	$0 * P(\text{NNP}   t_{\text{Janet}}) * s(t_{\text{Janet}}   \text{Janet})$			
MD	0	...			
VB	0	...			
JJ	0	...			
NN	0	...			
RB	0	...			
DT	0	...			


# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 	0			
MD	0	$P(\text{will}   \text{MD}) * P(\text{MD}   t_{\text{Janet}}) * s(t_{\text{Janet}}   \text{Janet})$			
VB	0	...			
JJ	0	...			
NN	0	...			
RB	0	...			
DT	0	...			

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 	0			
MD	0	3.004e-8			
VB	0	...			
JJ	0	...			
NN	0	...			
RB	0	...			
DT	0	...			

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 	0			
MD	0	3.004e-8			
VB	0	2.231e-13			
JJ	0	0			
NN	0	1.034e-10			
RB	0	0			
DT	0	0			

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0			
MD	0	3.004e-8			
VB	0	2.231e-13			
JJ	0	0			
NN	0	1.034e-10			
RB	0	0			
DT	0	0			

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0		
MD	0	3.004e-8	0		
VB	0	2.231e-13	$P(\text{back} \text{VB}) * P(\text{VB} t_{\text{will}}) * s(t_{\text{will}} \text{will})$		
JJ	0	0			
NN	0	1.034e-10			
RB	0	0			
DT	0	0			

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0		
MD	0	3.004e-8	0		
VB	0	2.231e-13	MD: 1.6e-11 VB: 7.5e-19 NN: 9.7e-17		
JJ	0	0			
NN	0	1.034e-10			
RB	0	0			
DT	0	0			



# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0		
MD	0	3.004e-8	0		
VB	0	2.231e-13	<b>MD: 1.6e-11</b> <b>VB: 7.5e-19</b> <b>NN: 9.7e-17</b>		
JJ	0	0			
NN	0	1.034e-10			
RB	0	0			
DT	0	0			

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0		
MD	0	3.004e-8	0		
VB	0	2.231e-13	1.6e-11		
JJ	0	0			
NN	0	1.034e-10			
RB	0	0			
DT	0	0			

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0		
MD	0	3.004e-8	0		
VB	0	2.231e-13	1.6e-11		
JJ	0	0	5.1e-15		
NN	0	1.034e-10	5.4e-15		
RB	0	0	5.3e-11		
DT	0	0	0		

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0	2.5e-17	
MD	0	3.004e-8	0	0	
VB	0	2.231e-13	1.6e-11	0	
JJ	0	0	5.1e-15	5.2e-16	
NN	0	1.034e-10	5.4e-15	5.9e-18	
RB	0	0	5.3e-11	0	
DT	0	0	0	1.8e-12	

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0	2.5e-17	0
MD	0	3.004e-8	0	0	0
VB	0	2.231e-13	1.6e-11	0	1.0e-20
JJ	0	0	5.1e-15	5.2e-16	0
NN	0	1.034e-10	5.4e-15	5.9e-18	2.0e-15
RB	0	0	5.3e-11	0	0
DT	0	0	0	1.8e-12	0

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0	2.5e-17	0
MD	0	3.004e-8	0	0	0
VB	0	2.231e-13	1.6e-11	0	1.0e-20
JJ	0	0	5.1e-15	5.2e-16	0
NN	0	1.034e-10	5.4e-15	5.9e-18	<b>2.0e-15</b>
RB	0	0	5.3e-11	0	0
DT	0	0	0	1.8e-12	0

# The viterbi algorithm

	Janet	will	back	the	bill
NNP	8.8544e-06 ●	0	0	2.5e-17	0
MD	0	3.004e-8	0	0	0
VB	0	2.231e-13	1.6e-11	0	1.0e-20
JJ	0	0	5.1e-15	5.2e-16	0
NN	0	1.034e-10	5.4e-15	5.9e-18	<b>2.0e-15</b>
RB	0	0	5.3e-11	0	0
DT	0	0	0	1.8e-12	0

# The viterbi algorithm

	Janet/ <b>NNP</b>	will/ <b>MD</b>	back/ <b>VB</b>	the/ <b>DT</b>	bill/ <b>NN</b>
NNP	8.8544e-06 ●	0	0	2.5e-17	0
MD	0	3.004e-8	0	0	0
VB	0	2.231e-13	1.6e-11	0	1.0e-20
JJ	0	0	5.1e-15	5.2e-16	0
NN	0	1.034e-10	5.4e-15	5.9e-18	<b>2.0e-15</b>
RB	0	0	5.3e-11	0	0
DT	0	0	0	1.8e-12	0



# The Viterbi algorithm

- Complexity:  $O(T^2N)$ , where  $T$  is the size of the tagset and  $N$  is the length of the sequence.
  - \*  $T * N$  matrix, each cell performs  $T$  operations.
- Why does it work?
  - \* Because of the **independence assumptions** that decompose the problem (specifically, the Markov property). Without these, we cannot apply DP.

# Viterbi Pseudocode

```
alpha = np.zeros(M, T)
for t in range(T):
    alpha[1, t] = pi[t] * O[w[1], t]

for i in range(2, M):
    for t_i in range(T):
        for t_last in range(T):           # t_last means t_{i-1}
            s = alpha[i-1, t_last] * A[t_last, t_i] * O[w[i], t_i]
            if s > alpha[i, t_i]:
                alpha[i, t_i] = s
                back[i, t_i] = t_last

best = np.max(alpha[M-1, :])
return backtrack(best, back)
```

- Good practice: work with **log** probabilities to prevent underflow (multiplications become sums)
- Vectorisation (use matrix-vector operations)

# HMMs in practice

- We saw HMM taggers based on **bigrams**. State-of-the-art use tag **trigrams**.
  - \*  $P(\mathbf{t}) = \prod_{i=1}^n P(t_i | t_{i-1}, t_{i-2})$  Viterbi now  $O(T^3N)$
- Need to deal with sparsity: some tag trigram sequences might not be present in training data
  - \* Backoff:  $P(t_i | t_{i-1}, t_{i-2}) = \lambda_3 \hat{P}(t_i | t_{i-1}, t_{i-2}) + \lambda_2 \hat{P}(t_i | t_{i-1}) + \lambda_1 \hat{P}(t_i)$
  - \*  $\lambda_1 + \lambda_2 + \lambda_3 = 1$
  - \* Can learn the weights using **deleted interpolation**.
- With additional features, reach 96.7% accuracy on Penn Treebank (Brants, 2000)

# Other variant Taggers

- HMM is **generative**,  $P(\mathbf{t}, \mathbf{w})$ , 'creates' the input
  - allows for unsupervised HMMs: learn model without any tagged data!
- **Discriminative** models also popular, modelling  $P(\mathbf{t} / \mathbf{w})$  directly
  - supports richer feature set, generally better accuracy when trained over large supervised datasets
  - E.g., Maximum Entropy Markov Model (MEMM), Conditional random field (CRF), Connectionist Temporal Classification (CTC)
  - Most *deep learning* models of sequences are discriminative (e.g., encoder-decoders for translation), similar to an MEMM

# HMMs in NLP

- HMMs are highly effective for part-of-speech tagging
  - trigram HMM gets 96.5% accuracy (TnT)
  - related models are state of the art
    - MEMMs            97%
    - CRFs             97.6%
    - Deep CRF        97.9%
  - *English Penn Treebank* tagging accuracy  
[https://aclweb.org/aclwiki/index.php?title=POS\\_Tagging\\_\(State\\_of\\_the\\_art\)](https://aclweb.org/aclwiki/index.php?title=POS_Tagging_(State_of_the_art))
- Other sequence labelling tasks
  - named entity recognition, shallow parsing, alignment ...
  - In other fields: DNA, protein sequences, image lattices...

# A final word

- HMMs are a simple, yet effective way to perform sequence labelling.
- Can still be competitive, and fast. Natural baseline for other sequence labelling tasks.
- Main drawback: not very flexible in terms of feature representation, compared to MEMMs and CRFs.

# Readings

- JM3 Appendix A A.1-A.2, A.4
- See also E18, parts of Chapter 7
- References:
  - \* Rabiner's HMM tutorial <http://tinyurl.com/2hqaf8>
  - \* Lafferty et al, Conditional random fields: Probabilistic models for segmenting and labeling sequence data (2001)