

Sequence Labelling

Dr. Uzair Ahmad

Program

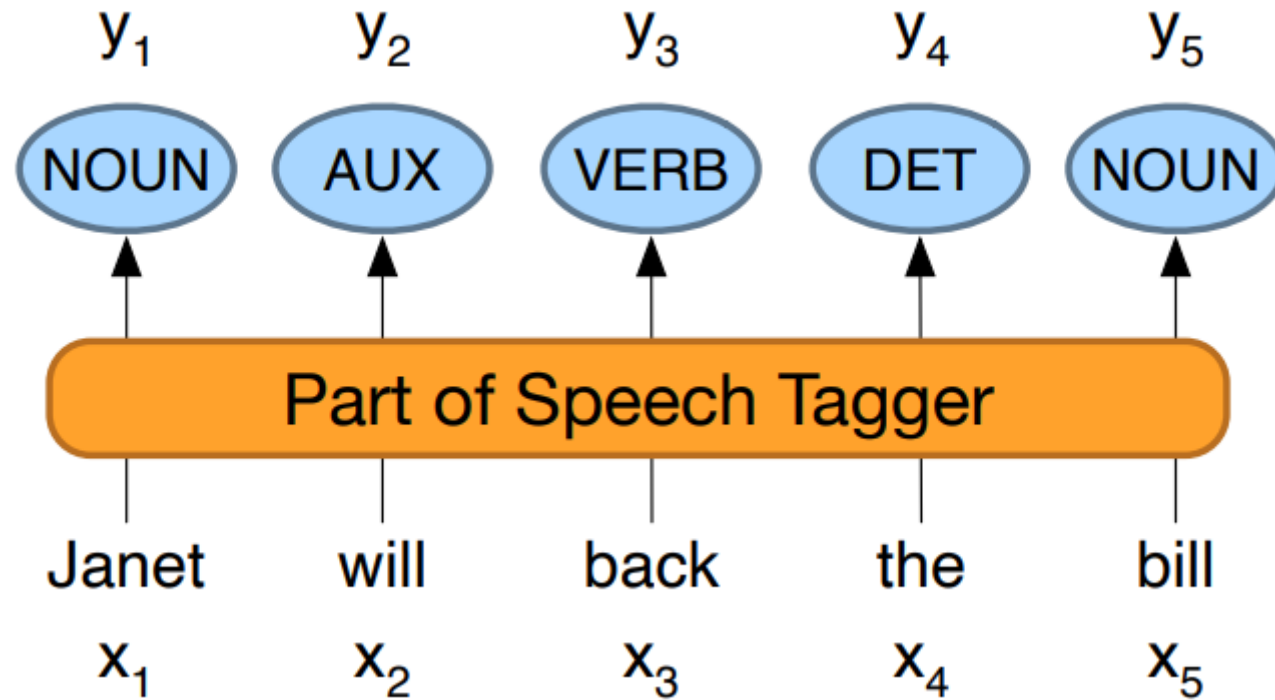
- Sequence labeling:
 - Part-of-speech tagging
 - NER
- Learning approaches
 - Classification approach
 - POS as structure prediction
 - Hidden Markov Models
 - The Viterbi algorithm
 - Discriminative approaches
 - RNNs and POS tagging

POS Tags

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	“to”	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>’s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past partici- ple	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one’s</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &</i>	WRB	wh-adverb	<i>how, where</i>

Penn Treebank part-of-speech tags.

POS Tagging Task



Sequence labeling

- Why difficult ?
 - They can fish. \rightarrow (N N V) (N V N)
 - Can of fish \rightarrow (N P N)

Sequence labeling

	They	can	fish
Possible assignments	V	V	V
	V	V	N
	V	N	N
	V	N	V
	N	N	N
	N	V	N
	N	N	V
	N	V	V

Classification approach

- $\mathbf{w} = (w_1, w_2, w_3, \dots, w_M)$
- Feature function
 - $f((w, m), y)$

$$f((\mathbf{w} = \textit{they can fish}, m = 1), \mathbf{N}) = (\textit{they}, \mathbf{N})$$

$$f((\mathbf{w} = \textit{they can fish}, m = 2), \mathbf{V}) = (\textit{can}, \mathbf{V})$$

$$f((\mathbf{w} = \textit{they can fish}, m = 3), \mathbf{V}) = (\textit{fish}, \mathbf{V}).$$

Classification approach

- Grammatical ambiguity
 - They **can** **fish**. \rightarrow (N N V) (N V N)
 - The **can** of **fish** \rightarrow (D N P N)
- The tagger must rely on context

Classification approach

- Context and Grammatical ambiguity

$$\begin{aligned} f((\boldsymbol{w} = \textit{they can fish}, 1), \text{N}) = \{ & (w_m = \textit{they}, y_m = \text{N}), \\ & (w_{m-1} = \square, y_m = \text{N}), \\ & (w_{m+1} = \textit{can}, y_m = \text{N}) \} \end{aligned}$$

$$\begin{aligned} f((\boldsymbol{w} = \textit{they can fish}, 2), \text{V}) = \{ & (w_m = \textit{can}, y_m = \text{V}), \\ & (w_{m-1} = \textit{they}, y_m = \text{V}), \\ & (w_{m+1} = \textit{fish}, y_m = \text{V}) \} \end{aligned}$$

$$\begin{aligned} f((\boldsymbol{w} = \textit{they can fish}, 3), \text{V}) = \{ & (w_m = \textit{fish}, y_m = \text{V}), \\ & (w_{m-1} = \textit{can}, y_m = \text{V}), \\ & (w_{m+1} = \blacksquare, y_m = \text{V}) \}. \end{aligned}$$

Classification approach

	1	2	3	4	5
W	The	old	man	the	boat.
POS	Det	Adj	N	Det	N
POS	Det	N	V	Det	N

Structure prediction

- Model the joint probability distribution $P(w, y)$
 - often using probabilistic graphical models.
- Set of tokens $\mathbf{w} = (w_1, w_2, w_3, \dots, w_M)$
- Set of possible tags $Y(w) = Y^M = (y_1, y_2, y_3, \dots, y_M)$
 - $Y = \{N, V, D, \dots\}$

$$\hat{y} = \operatorname{argmax}_{y \in Y(w)} \Psi(w, y)$$

Structure prediction

- Restricting the scoring function

$$\Psi(w, y) = \sum_{m=1}^{M+1} \psi(w, y_m, y_{m-1}, m)$$

$$\psi(w_{1:M}, y_m, y_{m-1}, m) = \theta \cdot f(w, y_m, y_{m-1}, m)$$

Structure prediction

$$\begin{aligned} f(\boldsymbol{w} = \textit{they can fish}, \boldsymbol{y} = \text{N V V}) &= \sum_{m=1}^{M+1} f(\boldsymbol{w}, y_m, y_{m-1}, m) \\ &= f(\boldsymbol{w}, \text{N}, \diamond, 1) \quad (w_m = \textit{they}, y_m = \text{N}) + (y_m = \text{N}, y_{m-1} = \diamond) \\ &+ f(\boldsymbol{w}, \text{V}, \text{N}, 2) \quad (w_m = \textit{can}, y_m = \text{V}) + (y_m = \text{V}, y_{m-1} = \text{N}) \\ &+ f(\boldsymbol{w}, \text{V}, \text{V}, 3) \quad (w_m = \textit{fish}, y_m = \text{V}) + (y_m = \text{V}, y_{m-1} = \text{V}) \\ &+ f(\boldsymbol{w}, \blacklozenge, \text{V}, 4) \quad (y_m = \blacklozenge, y_{m-1} = \text{V}) \end{aligned}$$

Structure prediction

- Inference by restricted scoring function

$$\begin{aligned}\Psi(w, y) &= \sum_{m=1}^{M+1} \psi(w, y_m, y_{m-1}, m) \\ &= \operatorname{argmax}_{y_{1:M}} \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) \\ &= \operatorname{argmax}_{y_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1}),\end{aligned}$$

Hidden Markov Models

- Probabilistic estimation of scores $s_m(\mathbf{y}, \mathbf{y}')$
- Naïve bayes classifier
 - $p(y|x) \propto p(x, y)$
- Probabilistic sequence labeling
 - $p(y|w) \propto p(y, w)$

$$\Psi(w, y) = \sum_{m=1}^{M+1} \psi(w, y_m, y_{m-1}, m)$$

Hidden Markov Models

- Probabilistic sequence labeling
 - $p(y|w) \propto p(y, w)$

$$p(y|w) = \prod_{m=1}^M p(w_m | y_m)$$

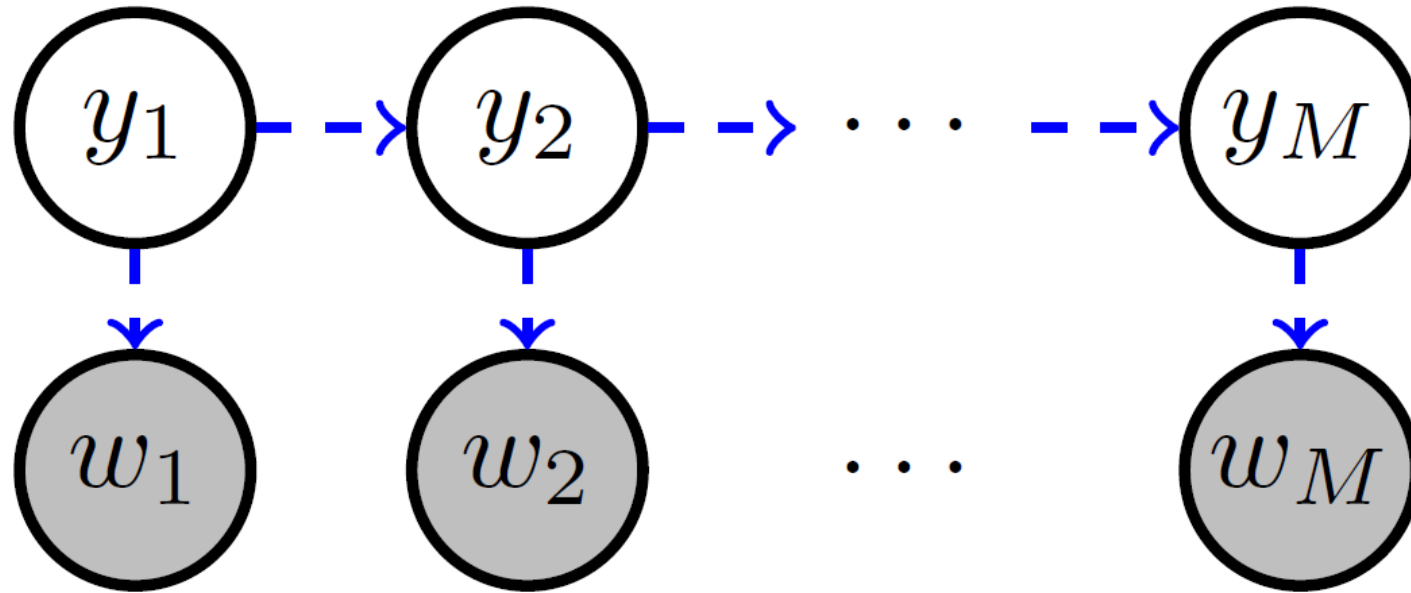
each token depends **only** on its tag

$$p(y) = \prod_{m=1}^M p(y_m | y_{m-1})$$

each tag depends **only** on its predecessor

Hidden Markov Models

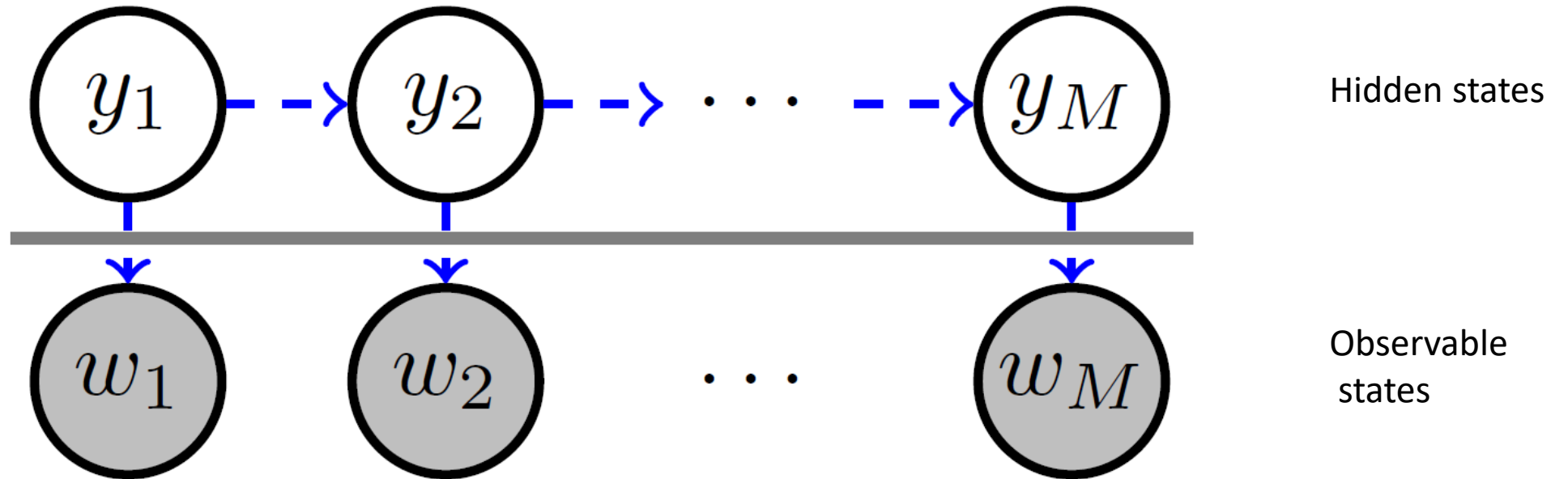
$$p(y) = \prod_{m=1}^M p(y_m | y_{m-1})$$



$$p(y|w) = \prod_{m=1}^M p(w_m | y_m)$$

Hidden Markov Models

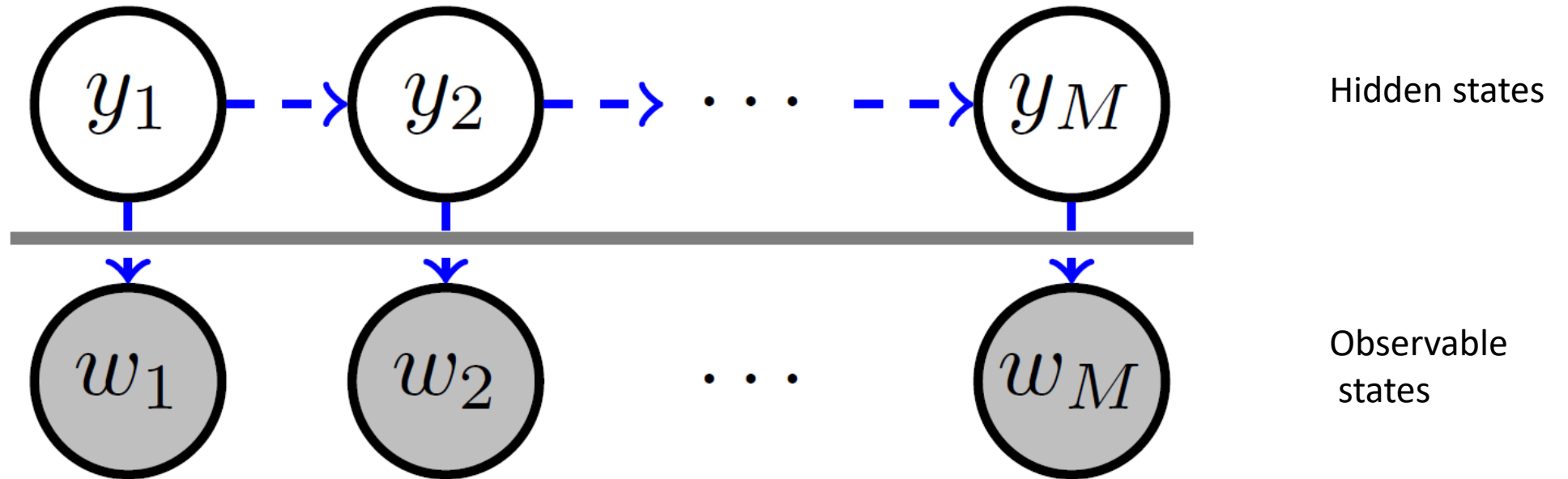
$$p(y) = \prod_{m=1}^M p(y_m | y_{m-1}) \rightarrow \text{Transition probabilities}$$



$$p(y|w) = \prod_{m=1}^M p(w_m | y_m) \rightarrow \text{Emission probabilities}$$

Hidden Markov Models

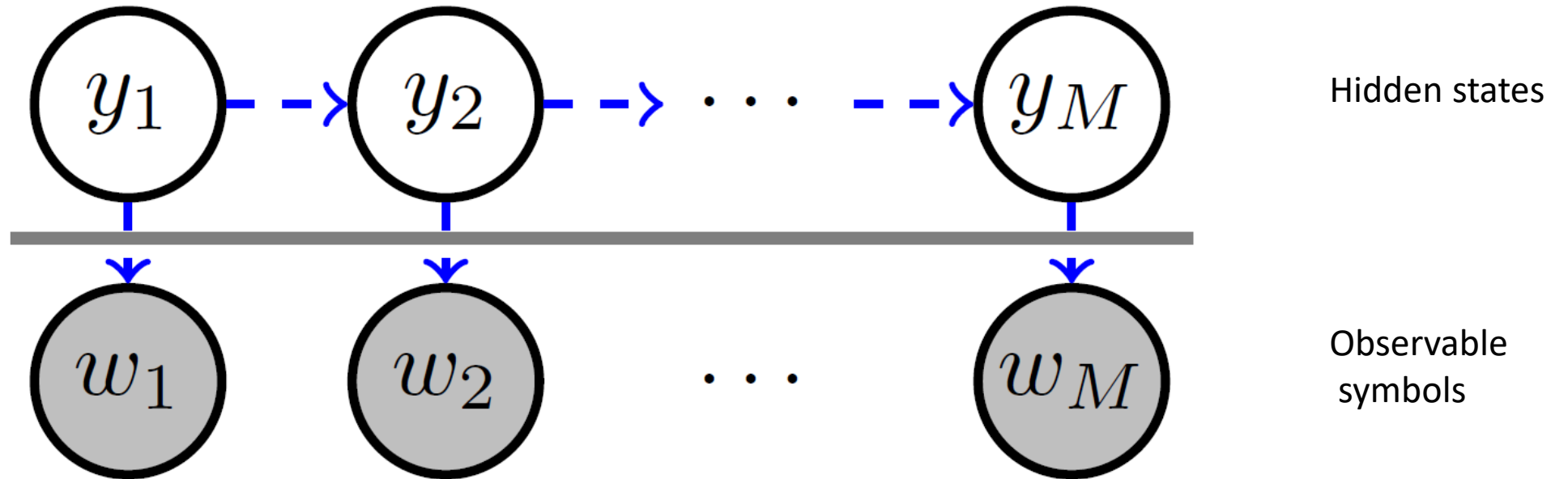
$$p_t(y) = \prod_{m=1}^M p(y_m | y_{m-1}; \lambda) \rightarrow \text{Transition probabilities}$$



$$p_e(y|w) = \prod_{m=1}^M p(w_m | y_m; \phi) \rightarrow \text{Emission probabilities}$$

Hidden Markov Models

$$p_t(y) = \prod_{m=1}^M p(y_m | y_{m-1}; \lambda) \rightarrow \text{Transition probabilities} \quad \lambda_{k,k'} \triangleq \Pr(Y_m = k' | Y_{m-1} = k) = \frac{\text{count}(Y_m = k', Y_{m-1} = k)}{\text{count}(Y_{m-1} = k)}$$



$$p_e(y|w) = \prod_{m=1}^M p(w_m | y_m; \phi) \rightarrow \text{Emission probabilities} \quad \phi_{k,i} \triangleq \Pr(W_m = i | Y_m = k) = \frac{\text{count}(W_m = i, Y_m = k)}{\text{count}(Y_m = k)}$$

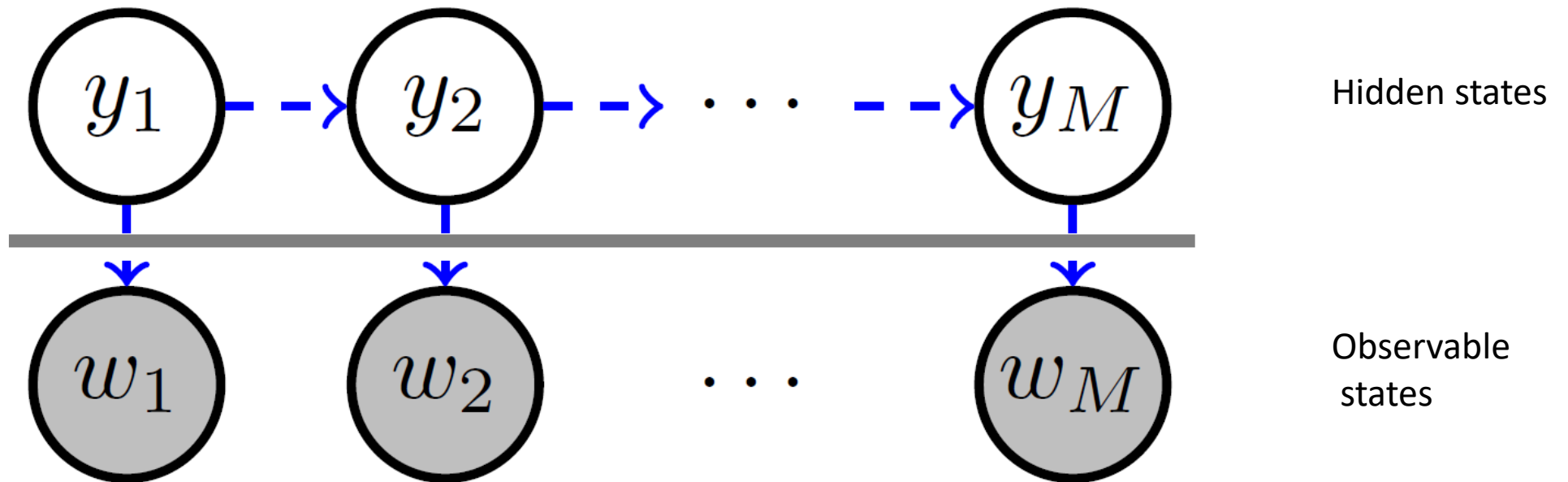
Hidden Markov Models

Inference

$$\hat{y} = \operatorname{argmax}_y p(y|w) \approx \max_y p(y, w)$$

$$p(y, w) = p(y | w) \times p(w) \propto p(y | w)$$

$$p(y, w) = p(w | y) \times p(y) \propto p(w | y)$$



Hidden Markov Models

Inference

$$\hat{y} = \operatorname{argmax}_y \log (p(y|w)) \approx \operatorname{max}_y \log (p(y, w))$$

$$\log (p(y, w)) = \log (p(y)) \times \log p(w|y)$$

$$\log (p(y, w)) = \sum_{m=1}^M \log p_Y(y_m | y_{m-1}) + \log p_{W|Y}(w_m | y_m)$$

$$\log (p(y, w)) = \sum_{m=1}^M \log \lambda_{y_m, y_{m-1}} + \log \phi_{y_m, w_m} \quad s_m(y_m, y_{m-1})$$

Viterbi Algorithm

$$\begin{aligned} v_m(y_m) &\triangleq \max_{\mathbf{y}_{1:m-1}} \sum_{n=1}^m s_n(y_n, y_{n-1}) \\ &= \max_{y_{m-1}} s_m(y_m, y_{m-1}) + \max_{\mathbf{y}_{1:m-2}} \sum_{n=1}^{m-1} s_n(y_n, y_{n-1}) \\ &= \max_{y_{m-1}} s_m(y_m, y_{m-1}) + v_{m-1}(y_{m-1}). \end{aligned}$$

Viterbi Algorithm

for $k \in \{0, \dots, K\}$ do $v_1(k) = s_1(k, \diamond)$	Initialize
for $m \in \{2, \dots, M\}$ do for $k \in \{0, \dots, K\}$ do $v_m(k) = \max_{k'} s_m(k, k') + v_{m-1}(k')$ $b_m(k) = \operatorname{argmax}_{k'} s_m(k, k') + v_{m-1}(k')$	Recursion
$y_M = \operatorname{argmax}_k s_{M+1}(\blacklozenge, k) + v_M(k)$	Termination
for $m \in \{M-1, \dots, 1\}$ do $y_m = b_m(y_{m+1})$ return $y_{1:M}$	Backtracking

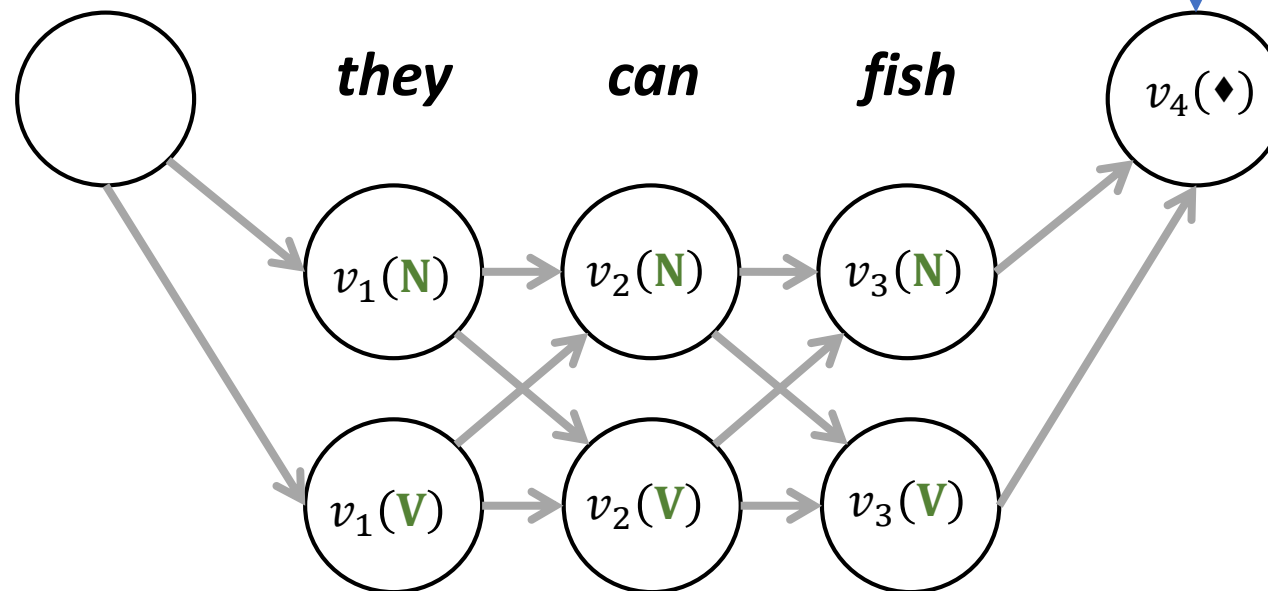
Viterbi Algorithm

	They	can	fish
Possible assignments	V	V	V
	V	V	N
	V	N	N
	V	N	V
	N	N	N
	N	V	N
	N	N	V
	N	V	V

Viterbi Algorithm

```
for  $k \in \{0, \dots, K\}$  do  
   $v_1(k) = s_1(k, \diamond)$ 
```

```
for  $m \in \{2, \dots, M\}$  do  
  for  $k \in \{0, \dots, K\}$  do  
     $v_m(k) = \max_{k'} s_m(k, k') + v_{m-1}(k')$   
     $b_m(k) = \operatorname{argmax}_{k'} s_m(k, k') + v_{m-1}(k')$ 
```



Viterbi Algorithm

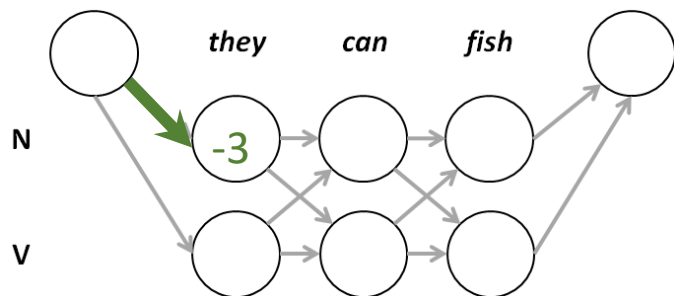
		<i>M</i>		
		<i>they</i>	<i>can</i>	<i>fish</i>
<i>K</i>	N	-2	-3	-3
	V	-10	-1	-3
		Emissions		

		<i>K - to</i>		
		N	V	◆
<i>K' - from</i>	◆	-1	-2	$-\infty$
	N	-3	-1	-1
	V	-1	-3	-1
		Transitions		

for $k \in \{0, \dots, K\}$ **do**
 $v_1(k) = s_1(k, \diamond)$

$$v_1(\mathbf{N}) = s_1(\mathbf{N}, \diamond) = -2 - 1 = -3$$

$$b_1(\mathbf{N}) = y_{m-1} = \diamond$$

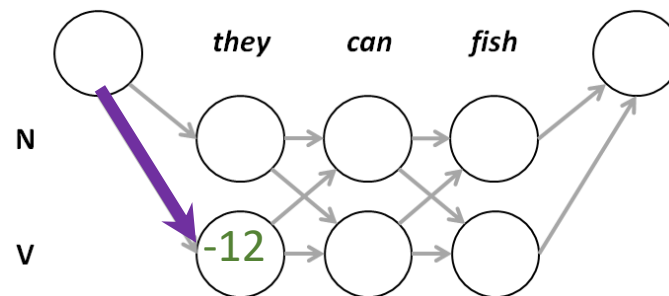


	<i>they</i>	<i>can</i>	<i>fish</i>
N	-2	-3	-3
V	-10	-1	-3

	N	V	◆
◇	-1	-2	$-\infty$
N	-3	-1	-1
V	-1	-3	-1

$$v_1(\mathbf{V}) = s_1(\mathbf{V}, \diamond) = -10 - 2 = -12$$

$$b_1(\mathbf{V}) = y_{m-1} = \diamond$$



	<i>they</i>	<i>can</i>	<i>fish</i>
N	-2	-3	-3
V	-10	-1	-3

	N	V	◆
◇	-1	-2	$-\infty$
N	-3	-1	-1
V	-1	-3	-1

```

for  $m \in \{2, \dots, M\}$  do
  for  $k \in \{0, \dots, K\}$  do
     $v_m(k) = \max_{k'} s_m(k, k') + v_{m-1}(k')$ 
     $b_m(k) = \operatorname{argmax}_{k'} s_m(k, k') + v_{m-1}(k')$ 

```

$$v_2(\mathbf{N}) = \max(v_1(\mathbf{N}) + s_2(\mathbf{N}, \mathbf{N}), v_1(\mathbf{V}) + s_2(\mathbf{N}, \mathbf{V}))$$

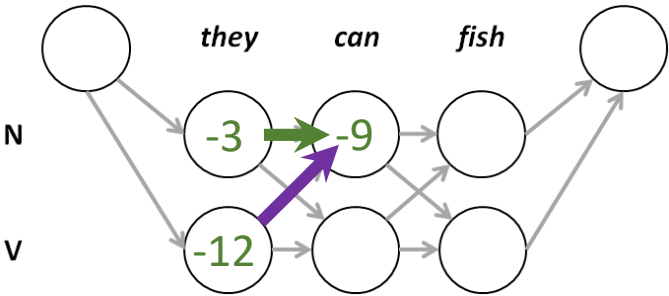
$$v_2(\mathbf{N}) = \max(-3 + s_2(\mathbf{N}, \mathbf{N}), -12 + s_2(\mathbf{N}, \mathbf{V}))$$

$$s_2(\mathbf{N}, \mathbf{N}) = -3 - 3$$

$$s_2(\mathbf{N}, \mathbf{V}) = -3 - 1$$

$$v_2(\mathbf{N}) = \max(-3 - 3 - 3, -12 - 3 - 1) = -9$$

$$b_2(\mathbf{N}) = y_{m-1} = \mathbf{N}$$



	they	can	fish
N	-2	-3	-3
V	-10	-1	-3

K - to

	N	V	◆
◇	-1	-2	$-\infty$
N	-3	-1	-1
V	-1	-3	-1

K' - from

```

for  $m \in \{2, \dots, M\}$  do
  for  $k \in \{0, \dots, K\}$  do
     $v_m(k) = \max_{k'} s_m(k, k') + v_{m-1}(k')$ 
     $b_m(k) = \operatorname{argmax}_{k'} s_m(k, k') + v_{m-1}(k')$ 

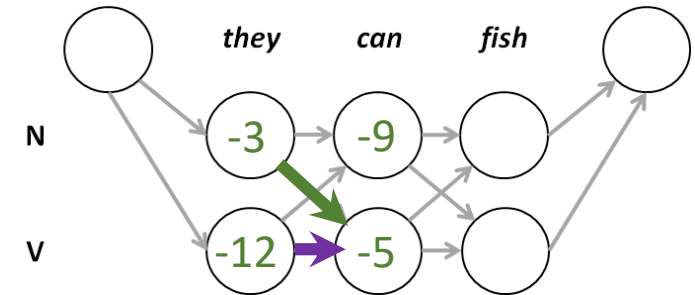
```

$$v_2(\mathbf{V}) = \max(v_1(\mathbf{N}) + s_2(\mathbf{V}, \mathbf{N}), v_1(\mathbf{V}) + s_2(\mathbf{V}, \mathbf{V}))$$

$$v_2(\mathbf{V}) = \max(-3 + \underset{-1-1}{s_2(\mathbf{V}, \mathbf{N})}, -12 + \underset{-1-3}{s_2(\mathbf{V}, \mathbf{V})})$$

$$v_2(\mathbf{V}) = \max(-3 \underset{\substack{\text{blue} \\ \swarrow}}{-1} \underset{\substack{\text{green} \\ \swarrow}}{-1}, -12 \underset{\substack{\text{blue} \\ \swarrow}}{-1} \underset{\substack{\text{purple} \\ \swarrow}}{-3}) = -5$$

$$b_2(\mathbf{V}) = y_{m-1} = \mathbf{N}$$



	<i>they</i>	<i>can</i>	<i>fish</i>
N	-2	-3	-3
V	-10	-1	-3

K - to

	N	V	◆
◇	-1	-2	$-\infty$
N	-3	-1	-1
V	-1	-3	-1

K' - from

```

for  $m \in \{2, \dots, M\}$  do
  for  $k \in \{0, \dots, K\}$  do
     $v_m(k) = \max_{k'} s_m(k, k') + v_{m-1}(k')$ 
     $b_m(k) = \operatorname{argmax}_{k'} s_m(k, k') + v_{m-1}(k')$ 

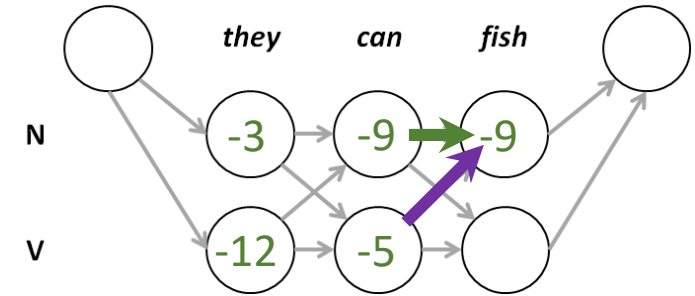
```

$$v_3(\mathbf{N}) = \max(v_2(\mathbf{N}) + s_3(\mathbf{N}, \mathbf{N}), v_2(\mathbf{V}) + s_3(\mathbf{N}, \mathbf{V}))$$

$$v_3(\mathbf{N}) = \max(-9 + \underset{-3-3}{s_3(\mathbf{N}, \mathbf{N})}, -5 + \underset{-3-1}{s_3(\mathbf{N}, \mathbf{V})})$$

$$v_3(\mathbf{N}) = \max(-9 \underset{-3}{-3} \underset{-3}{-3}, -5 \underset{-3}{-3} \underset{-1}{-1}) = -9$$

$$b_3(\mathbf{N}) = y_{m-1} = \mathbf{V}$$



	<i>they</i>	<i>can</i>	<i>fish</i>
N	-2	-3	-3
V	-10	-1	-3

K - to

	N	V	◆
◇	-1	-2	$-\infty$
N	-3	-1	-1
V	-1	-3	-1

K' - from

```

for  $m \in \{2, \dots, M\}$  do
  for  $k \in \{0, \dots, K\}$  do
     $v_m(k) = \max_{k'} s_m(k, k') + v_{m-1}(k')$ 
     $b_m(k) = \operatorname{argmax}_{k'} s_m(k, k') + v_{m-1}(k')$ 

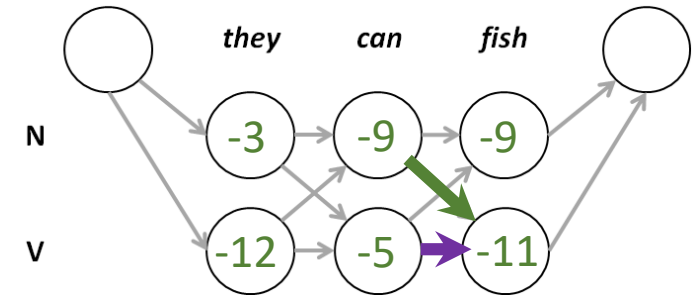
```

$$v_3(\mathbf{V}) = \max(v_2(\mathbf{N}) + s_3(\mathbf{V}, \mathbf{N}), v_2(\mathbf{V}) + s_3(\mathbf{V}, \mathbf{V}))$$

$$v_3(\mathbf{V}) = \max(-9 + \underset{-3-1}{s_3(\mathbf{V}, \mathbf{N})}, -5 + \underset{-3-3}{s_3(\mathbf{V}, \mathbf{V})})$$

$$v_3(\mathbf{V}) = \max(-9 \underset{-3}{-} \underset{-1}{-}, -5 \underset{-3}{-} \underset{-3}{-}) = -11$$

$$b_3(\mathbf{V}) = y_{m-1} = \mathbf{V}$$



	<i>they</i>	<i>can</i>	<i>fish</i>
N	-2	-3	-3
V	-10	-1	-3

K - to

	N	V	◆
◆	-1	-2	$-\infty$
N	-3	-1	-1
V	-1	-3	-1

K' - from

Termination

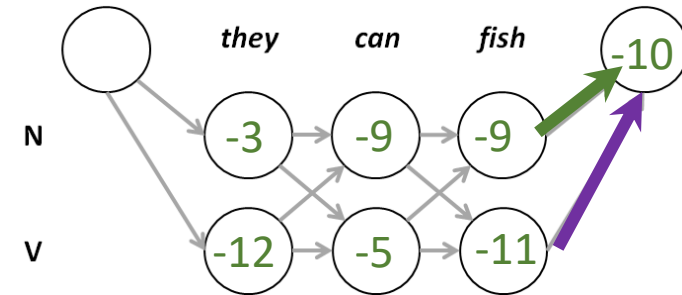
$$y_M = \operatorname{argmax}_k s_{M+1}(\blacklozenge, k) + v_M(k)$$

$$v_4(\blacklozenge) = \max(v_3(N) + s_4(\blacklozenge, N), v_3(V) + s_4(\blacklozenge, V))$$

$$v_4(\blacklozenge) = \max(-9 + \underset{0-1}{s_4(\blacklozenge, N)}, -11 + \underset{0-1}{s_4(\blacklozenge, V)})$$

$$v_4(\blacklozenge) = \max(-9 + 0-1, -11 + 0-1) = -10$$

$$b_4(\blacklozenge) = y_{m-1} = N$$



	they	can	fish
N	-2	-3	-3
V	-10	-1	-3

K' - from

	<i>K - to</i>		
	N	V	\blacklozenge
\lozenge	-1	-2	$-\infty$
N	-3	-1	-1
V	-1	-3	-1

Backtracking

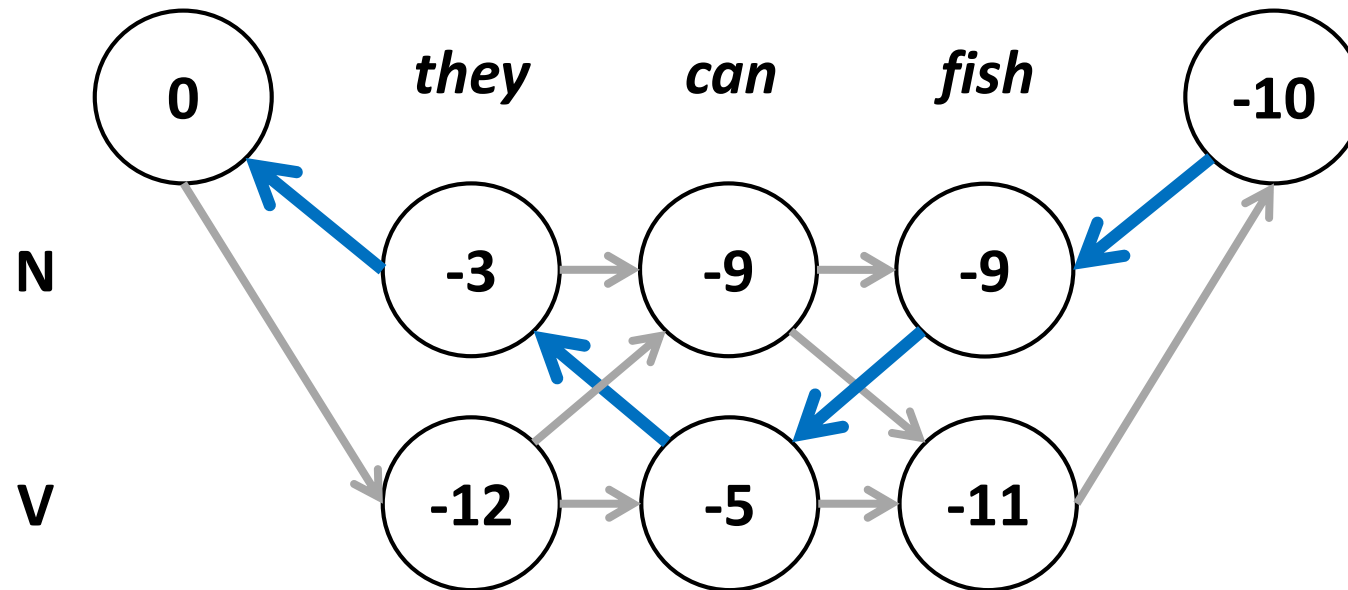
```
for  $m \in \{M - 1, \dots, 1\}$  do  
     $y_m = b_m(y_{m+1})$ 
```

$$y_0 = b_1(\mathbf{N}) = \diamond$$

$$y_2 = b_3(\mathbf{V}) = \mathbf{V}$$

$$y_1 = b_2(\mathbf{N}) = \mathbf{N}$$

$$y_3 = b_4(\diamond) = N$$



Viterbi Algorithm

- Restricting scoring function to local parts
 - Only pairs of adjacent tags
 - Akin to a bigram language model (over tags)
- Higher-order features ?

$$\Psi(w, y) = \sum_{m=1}^{M+2} \psi(w, y_m, y_{m-1}, y_{m-2}, m)$$

Discriminative approaches

- Structured Perceptron
 - Increase weights of correct classifications
 - Decrease weights of incorrect classification
- Structured SVM
 - Push class boundary away from training instances
- Conditional Random Fields
 - conditional probabilistic model for sequence labeling;
 - built on the logistic regression classifier

References

- Natural Language Processing | J. Eisenstein