

Diving into SyntaxNet

Kim Ho Yeob

NLP Procedure

Tokenizing

형태소 분석

POS-Tagging

품사 태깅

Syntactic Parsing

문장 구조 분석

Semantic Parsing

문장 의미 분석

SyntaxNet ?

Tokenizing

POS-Tagging

Syntactic Parsing

Semantic Parsing



SyntaxNet

SyntaxNet Main Paper

Globally Normalized Transition-Based Neural Networks

Globally Normalized

Tools for solving problems of local normalization

Transition-Based

Method of dependency parsing

Neural Networks

SyntaxNet Main Paper

Globally Normalized Transition-Based Neural Networks

1.

Globally Normalized

Tools for solving problems of local normalization

2.

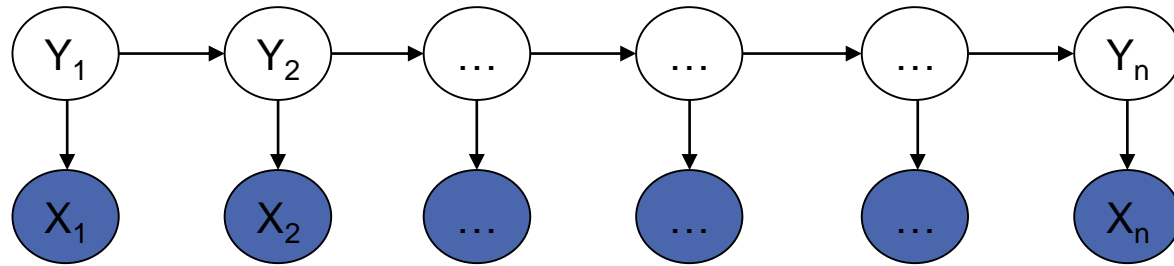
Transition-Based

Method of dependency parsing

Neural Networks

Globally Normalized?

- Globally Normalization starts from **Hidden Markov Model(HMM)**

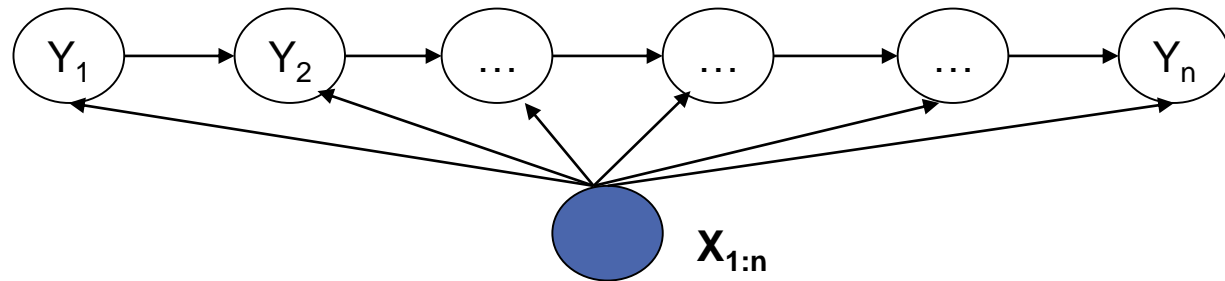


- Shortcomings of Hidden Markov Model
- HMM models direct dependence between each state and **only** its corresponding observation
 - NLP example: In a sentence segmentation task, segmentation may depend not just on a single word, but also on the features of the whole line such as line length, indentation, amount of white space, etc.

Globally Normalized?

- **Maximum Entropy Markov Model(MEMM)**

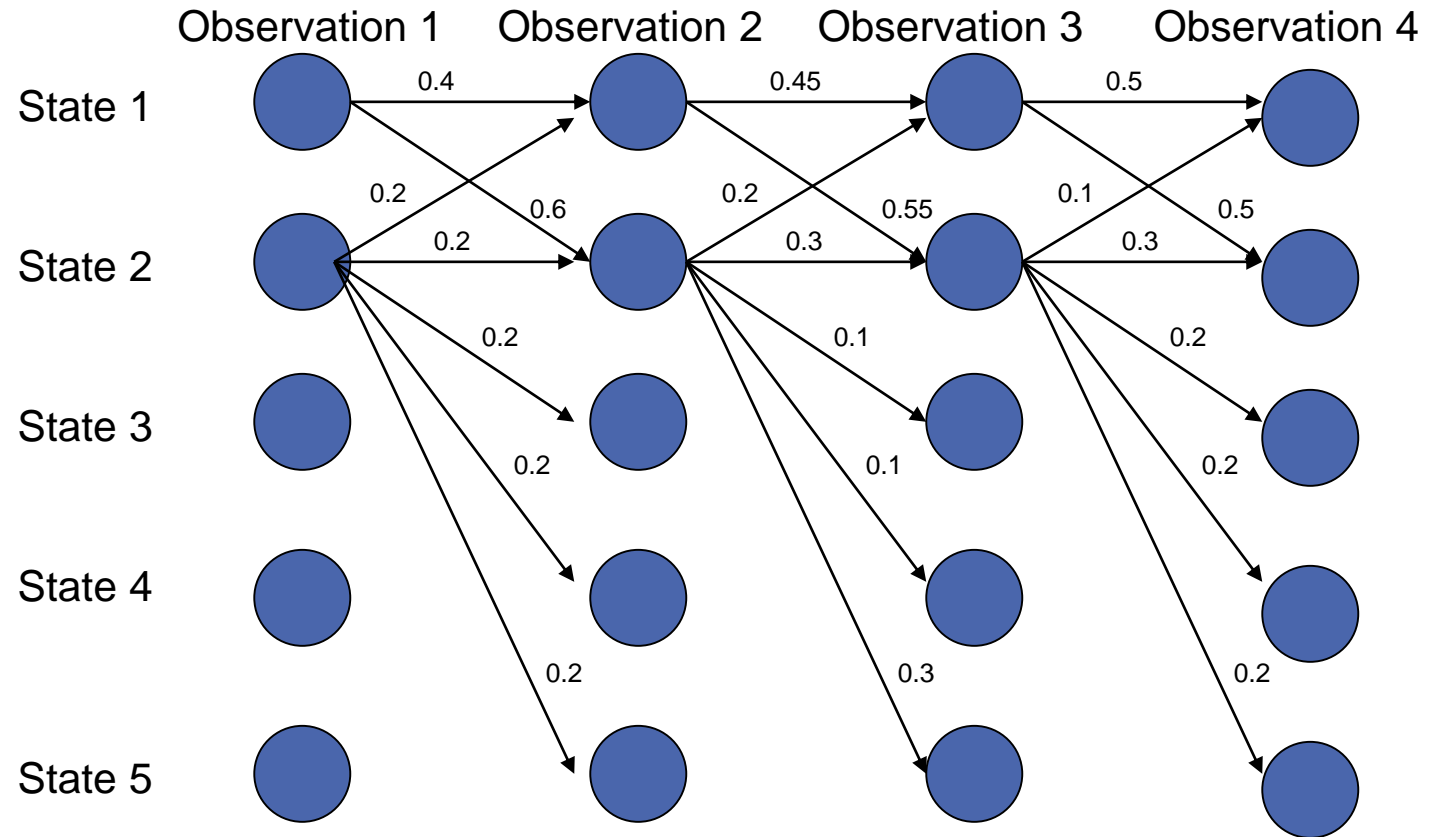
- Solution for shortcomings of HMM



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

- Models dependence between each state and the full observation sequence explicitly
 - More expressive than HMMs
- Discriminative model
 - Completely ignores modeling $P(\mathbf{X})$: saves modeling effort
 - Learning objective function consistent with predictive function: $P(\mathbf{Y}|\mathbf{X})$

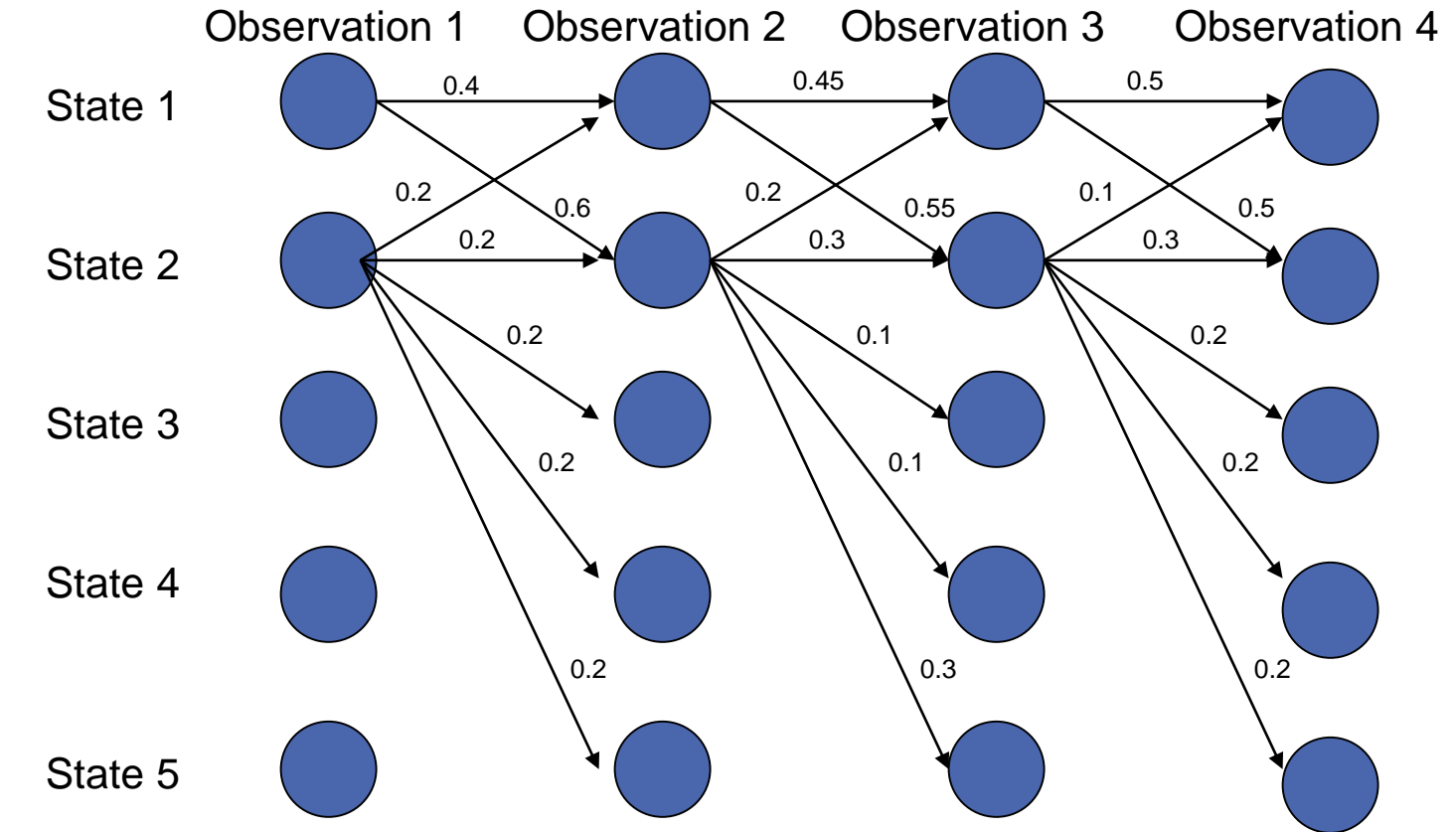
MEMM: Label bias problem



What the local transition probabilities say:

- State 1 almost always prefers to go to state 2
- State 2 almost always prefer to stay in state 2

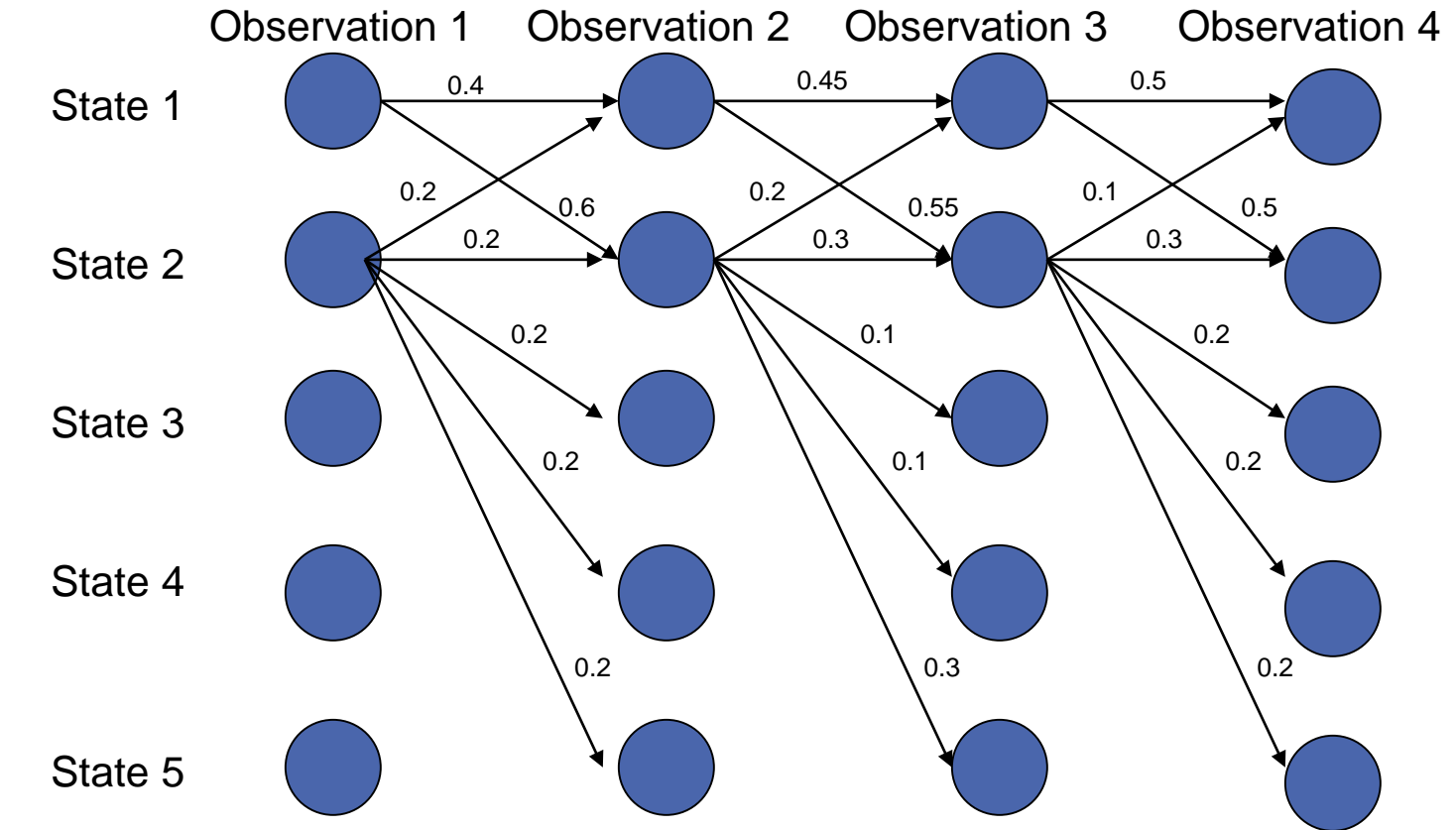
MEMM: Label bias problem



Probability of path 1-> 1-> 1-> 1:

- $0.4 \times 0.45 \times 0.5 = 0.09$

MEMM: Label bias problem



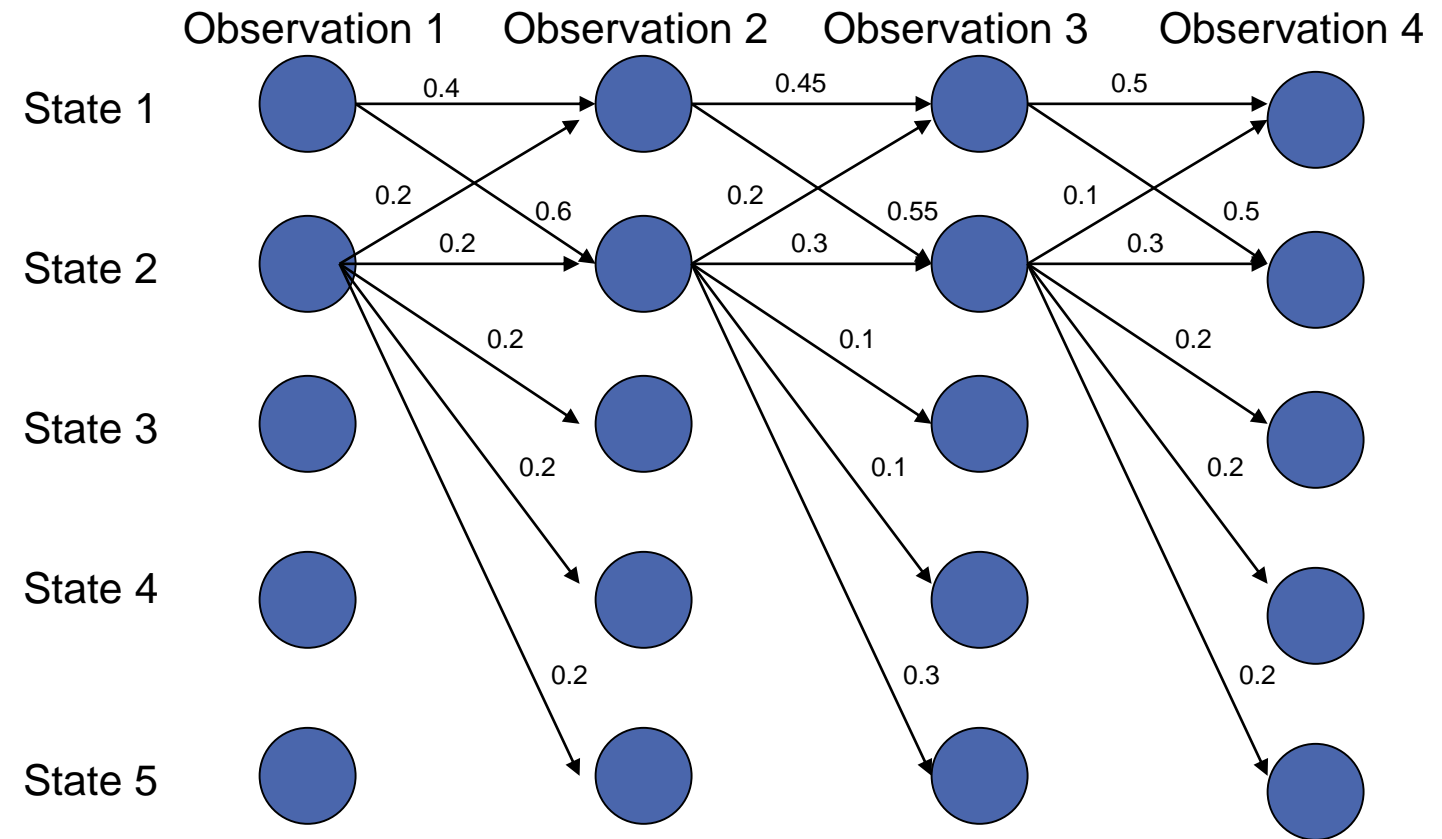
Probability of path 2->2->2->2 :

- $0.2 \times 0.3 \times 0.3 = 0.018$

Other paths:

1-> 1-> 1-> 1: 0.09

MEMM: Label bias problem



Probability of path 1->2->1->2:

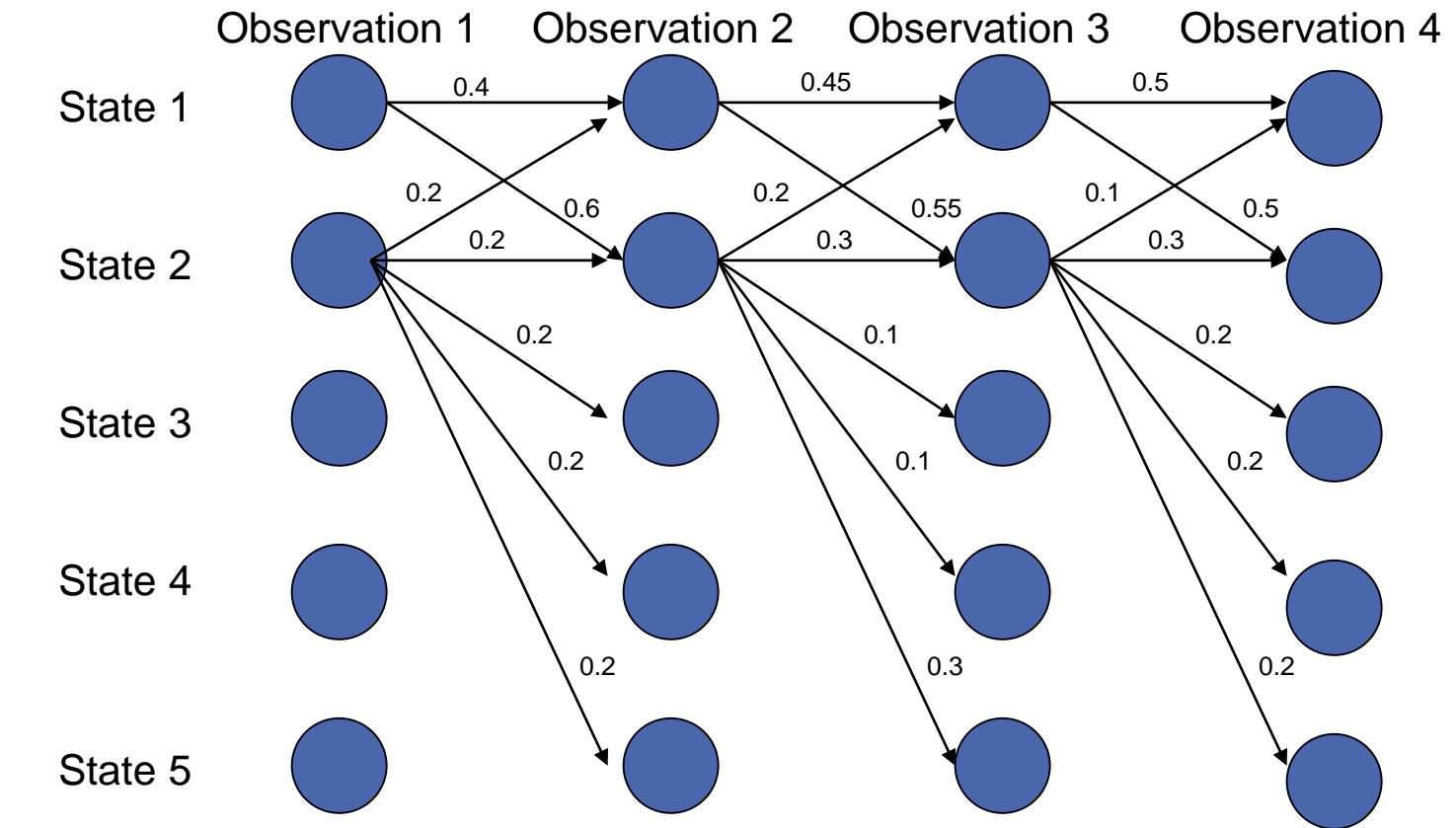
- $0.6 \times 0.2 \times 0.5 = 0.06$

Other paths:

1->1->1->1: 0.09

2->2->2->2: 0.018

MEMM: Label bias problem



Probability of path 1->1->2->2:

- $0.4 \times 0.55 \times 0.3 = 0.066$

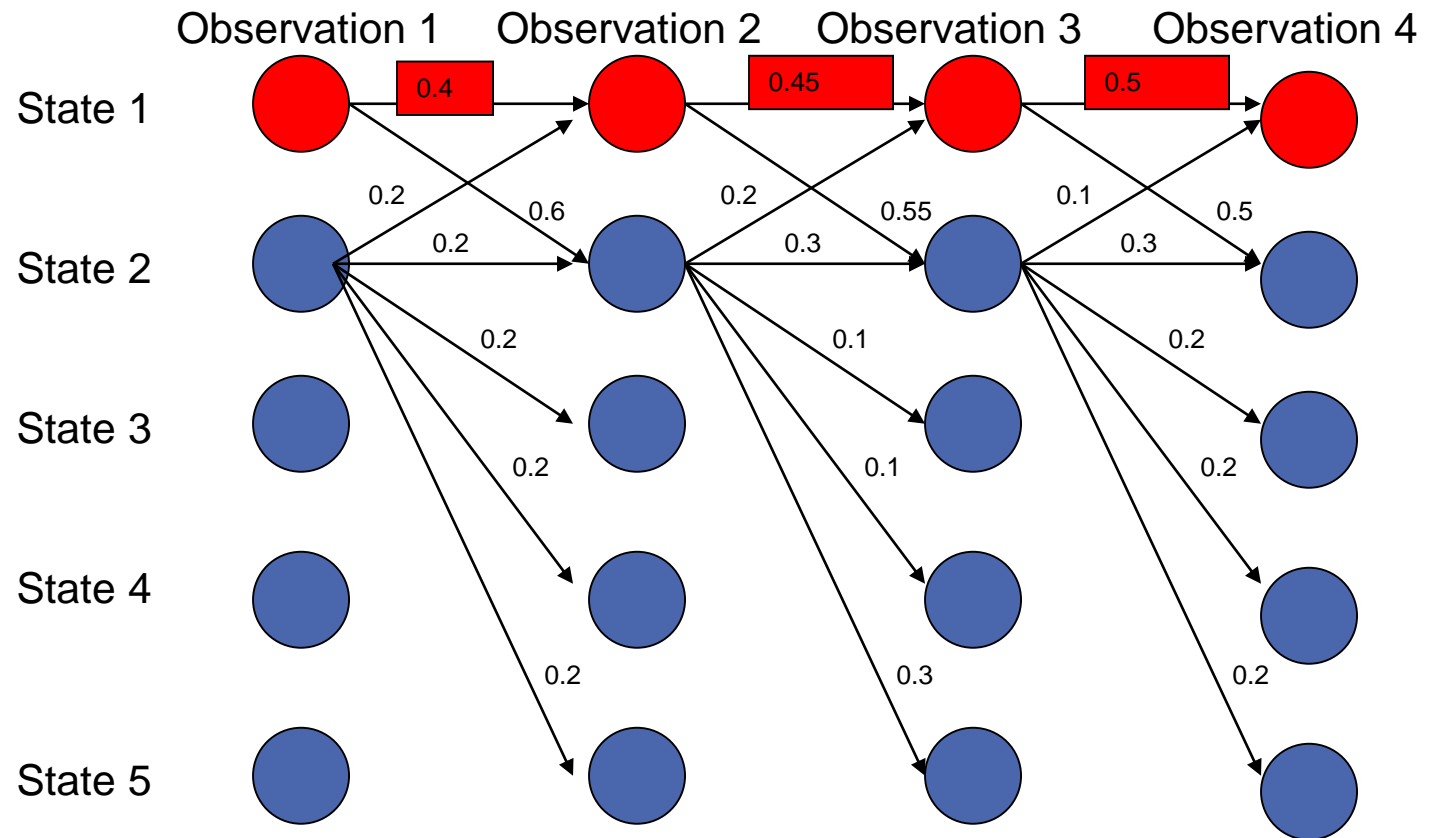
Other paths:

1->1->1->1: 0.09

2->2->2->2: 0.018

1->2->1->2: 0.06

MEMM: Label bias problem

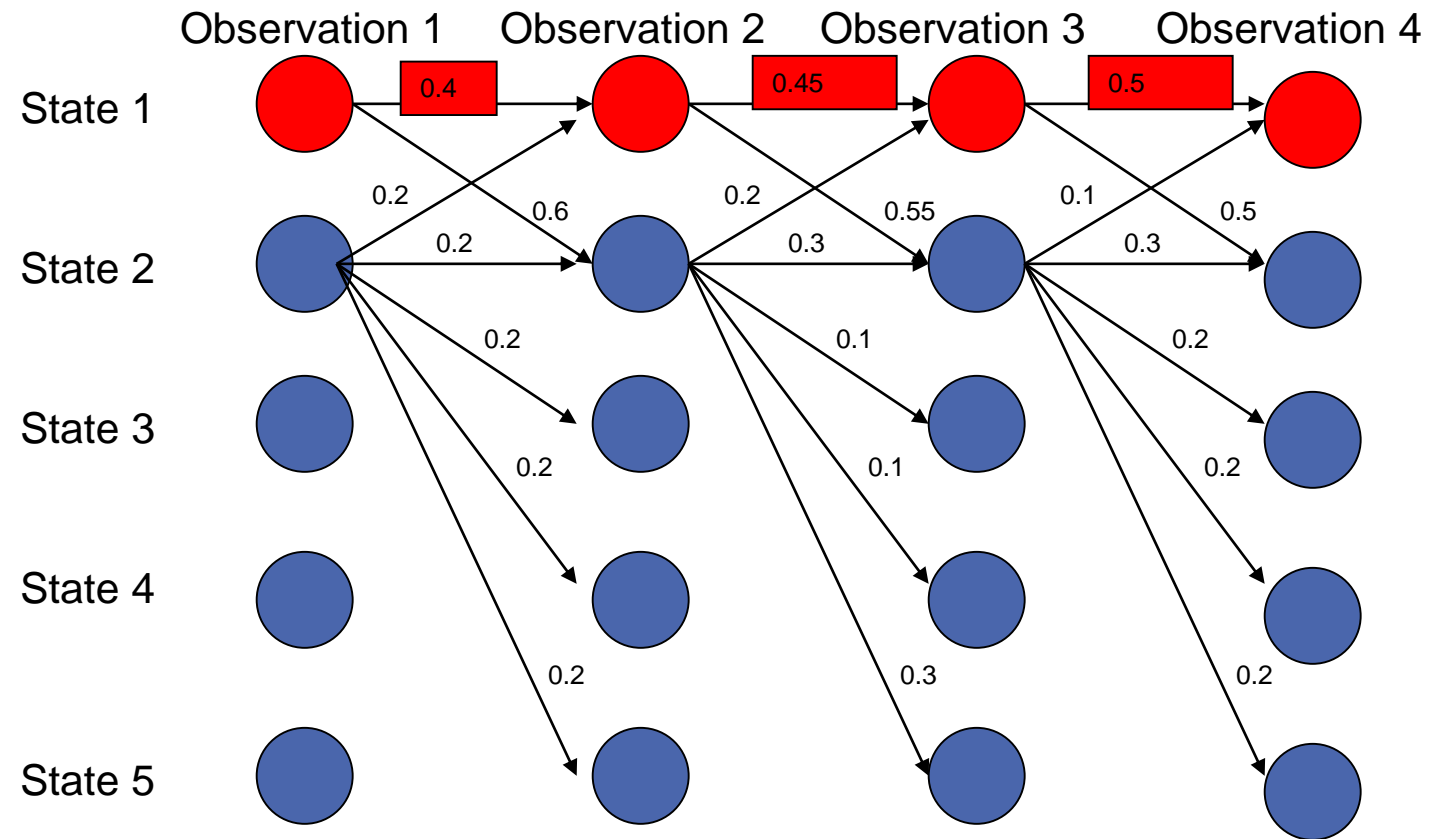


Most Likely Path: 1-> 1-> 1-> 1

- Although locally it seems state 1 wants to go to state 2 and state 2 wants to remain in state 2.

- **why?**

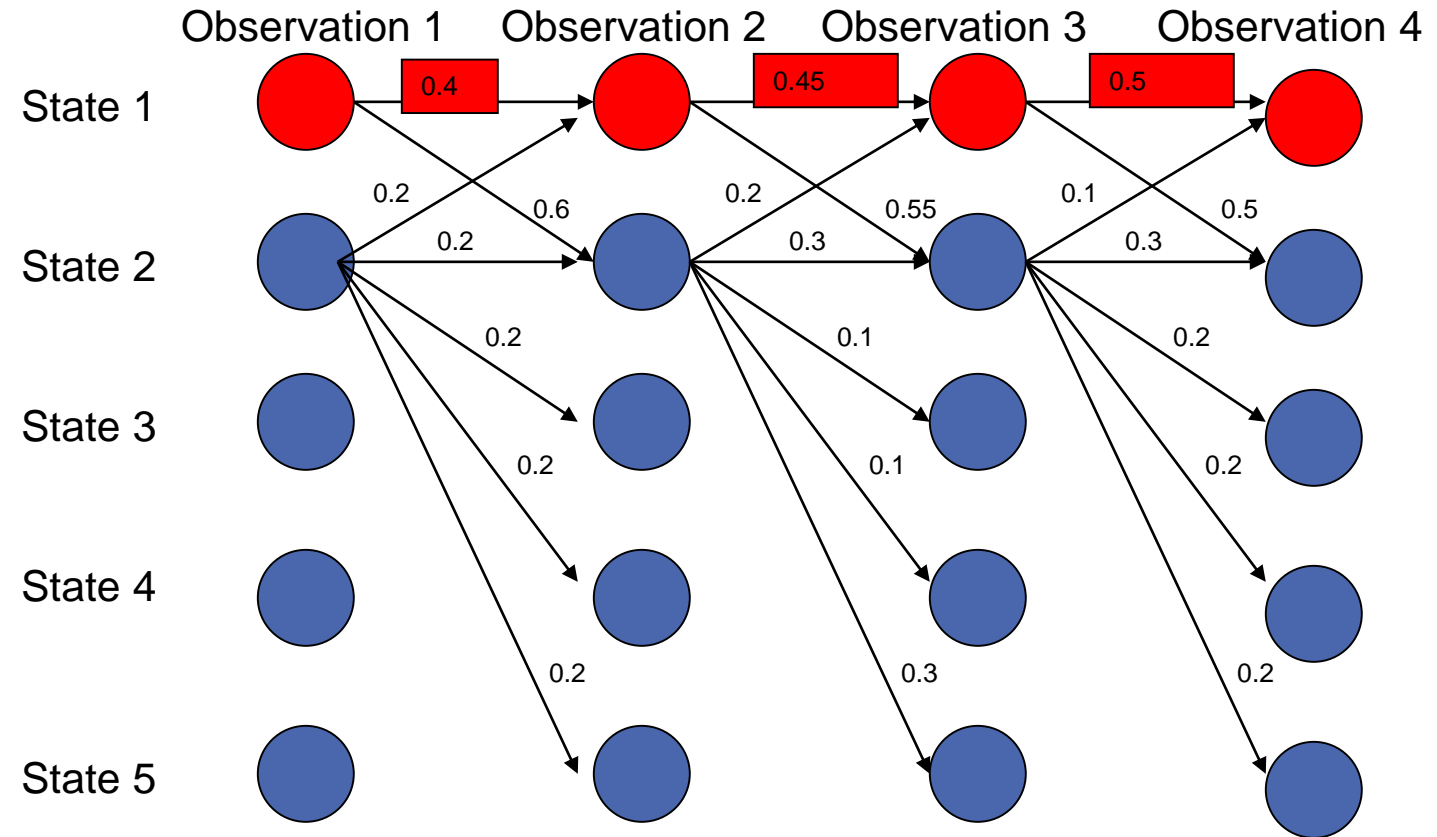
MEMM: Label bias problem



Most Likely Path: 1-> 1-> 1-> 1

- State 1 has only two transitions but state 2 has 5:
- Average transition probability from state 2 is lower

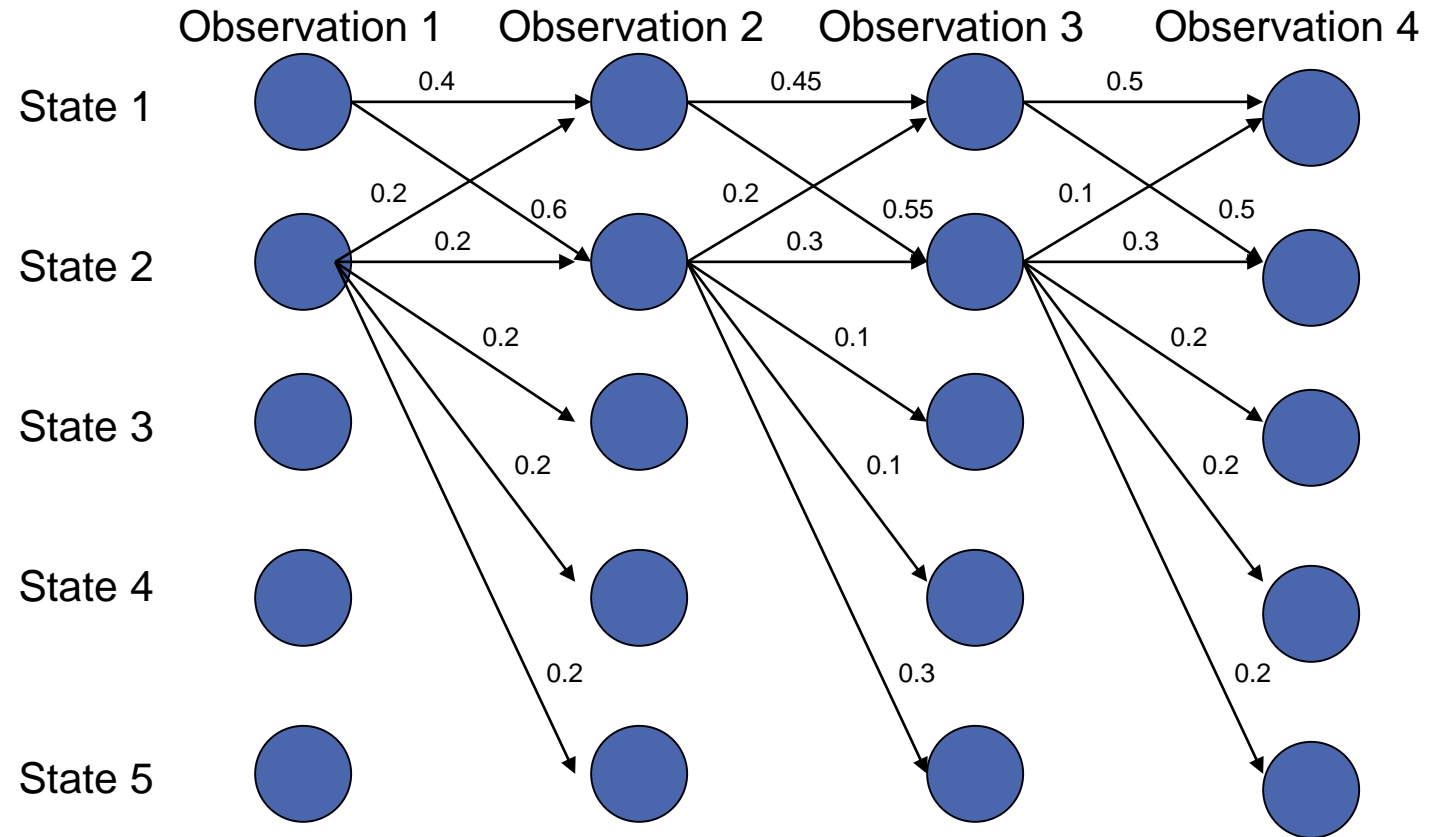
MEMM: Label bias problem



Label bias problem in MEMM:

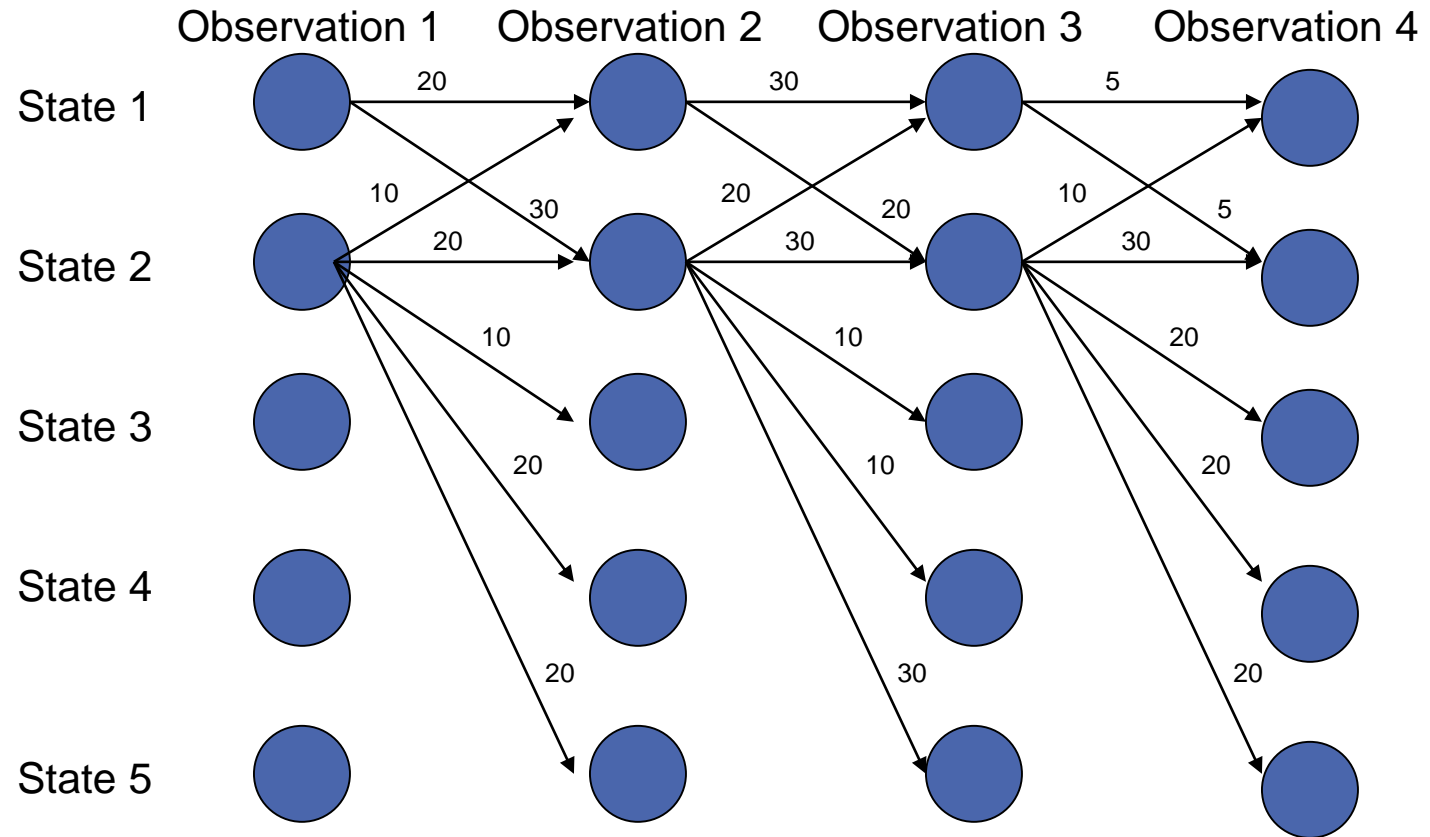
- Preference of states with lower number of transitions over others

Solution:
Do not
normalize
probabilities
locally



From local probabilities

Solution:
Do not
normalize
probabilities
locally

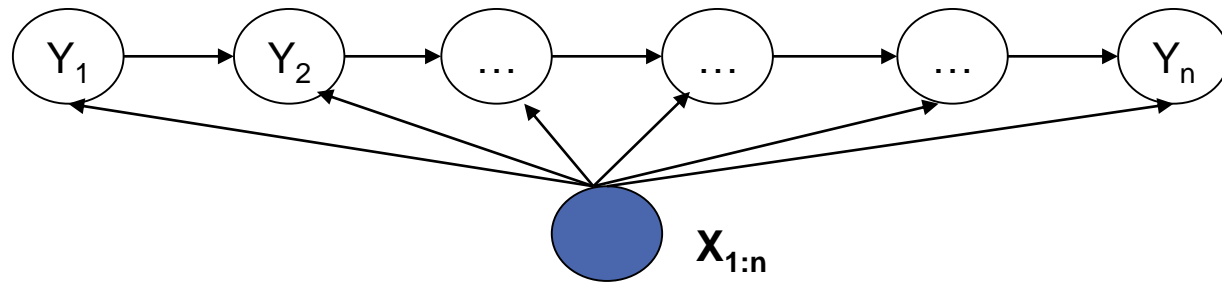


From local probabilities to local potentials

- States with lower transitions do not have an unfair advantage!

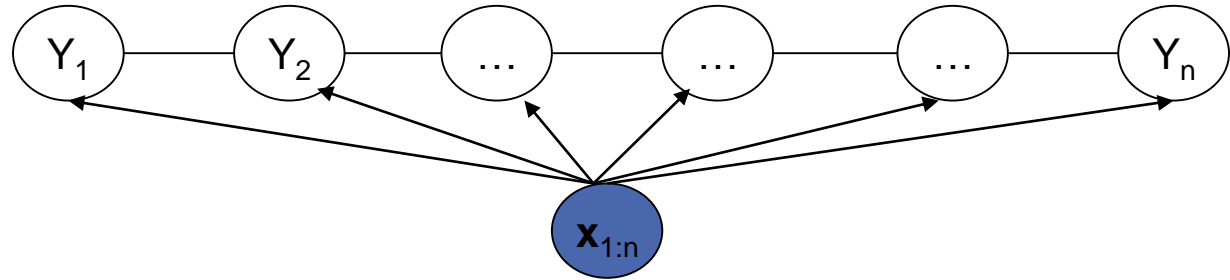
From MEMM

....



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

From MEMM to CRF

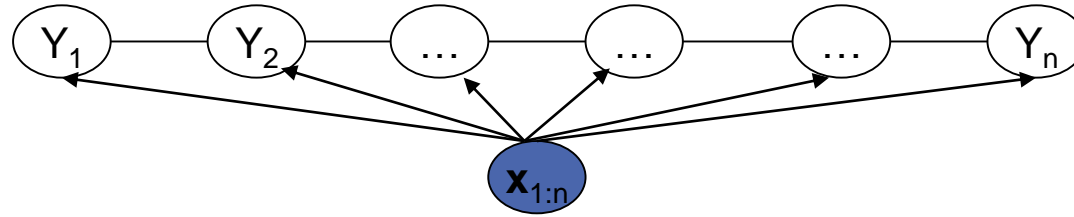


$$P(\mathbf{y}_{1:n}|\mathbf{x}_{1:n}) = \frac{1}{Z(\mathbf{x}_{1:n})} \prod_{i=1}^n \phi(y_i, y_{i-1}, \mathbf{x}_{1:n}) = \frac{1}{Z(\mathbf{x}_{1:n})} \prod_{i=1}^n \exp(\mathbf{w}^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}_{1:n}))$$

CRF is a partially directed model

- Discriminative model like MEMM
- Usage of global normalizer $Z(\mathbf{x})$ overcomes the label bias problem of MEMM
- Models the dependence between each state and the entire observation sequence (like MEMM)

Conditional Random Fields

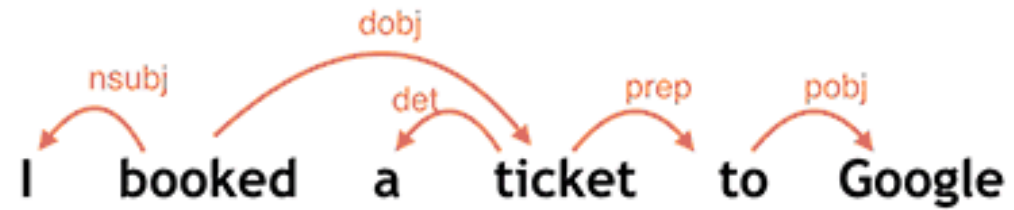


$$\begin{aligned} P(\mathbf{y}|\mathbf{x}) &= \frac{1}{Z(\mathbf{x})} \exp\left(\sum_{i=1}^n \left(\sum_k \lambda_k f_k(y_i, y_{i-1}, \mathbf{x}) + \sum_l \mu_l g_l(y_i, \mathbf{x})\right)\right) \\ &= \frac{1}{Z(\mathbf{x})} \exp\left(\sum_{i=1}^n (\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x}))\right) \end{aligned}$$

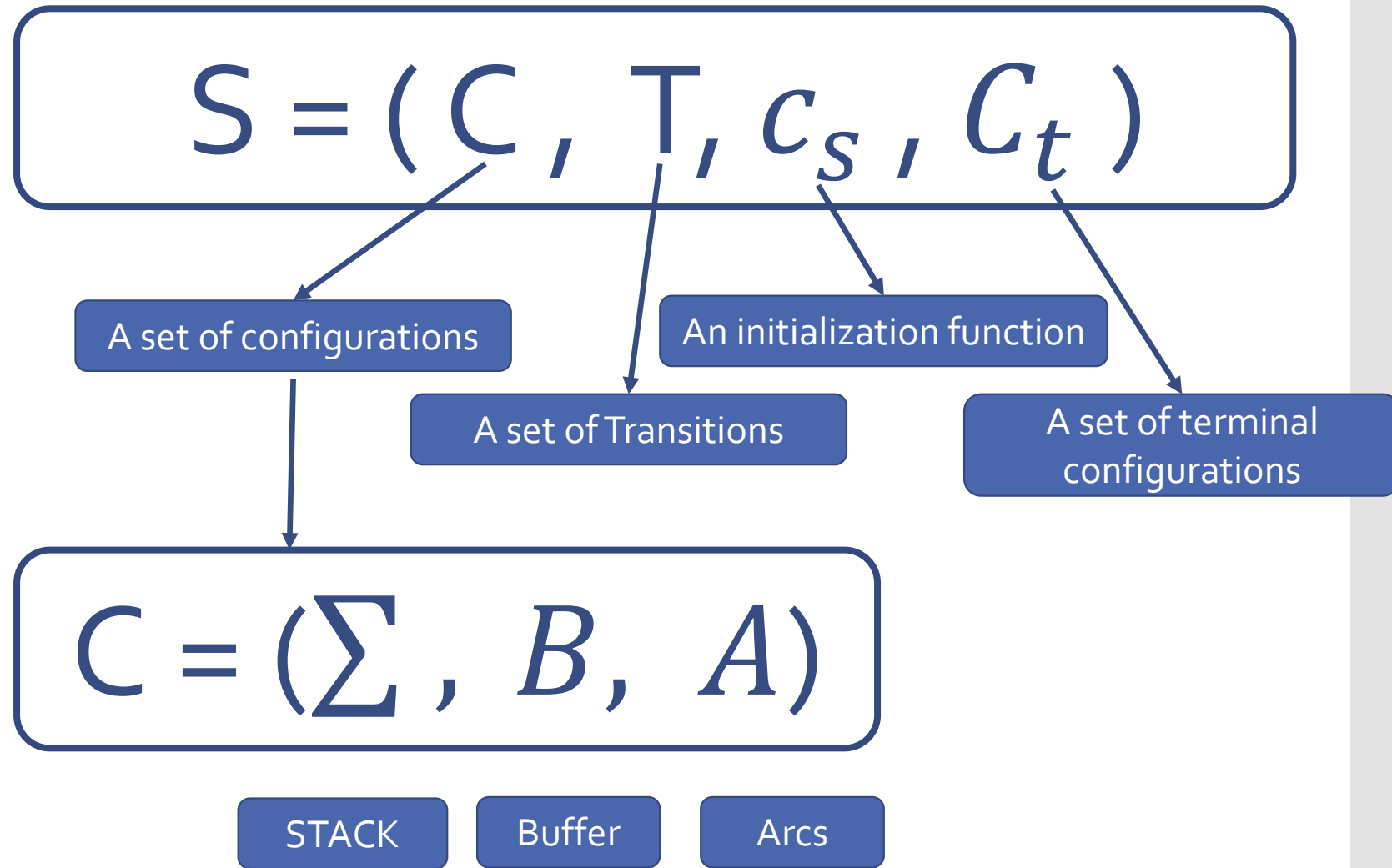
$$\text{where } Z(\mathbf{x}) = \sum_{\mathbf{y}} \exp\left(\sum_{i=1}^n (\lambda^T \mathbf{f}(y_i, y_{i-1}, \mathbf{x}) + \mu^T \mathbf{g}(y_i, \mathbf{x}))\right)$$

Dependency Parsing?

Dependency Parsing



Transition Systems



Transition Systems

Initialization: $c_s(x = x_1, \dots, x_n) = ([0], [1, \dots, n], \emptyset)$

Terminal: $C_t = \{c \in C \mid c = ([0], [], A)\}$

Transitions:

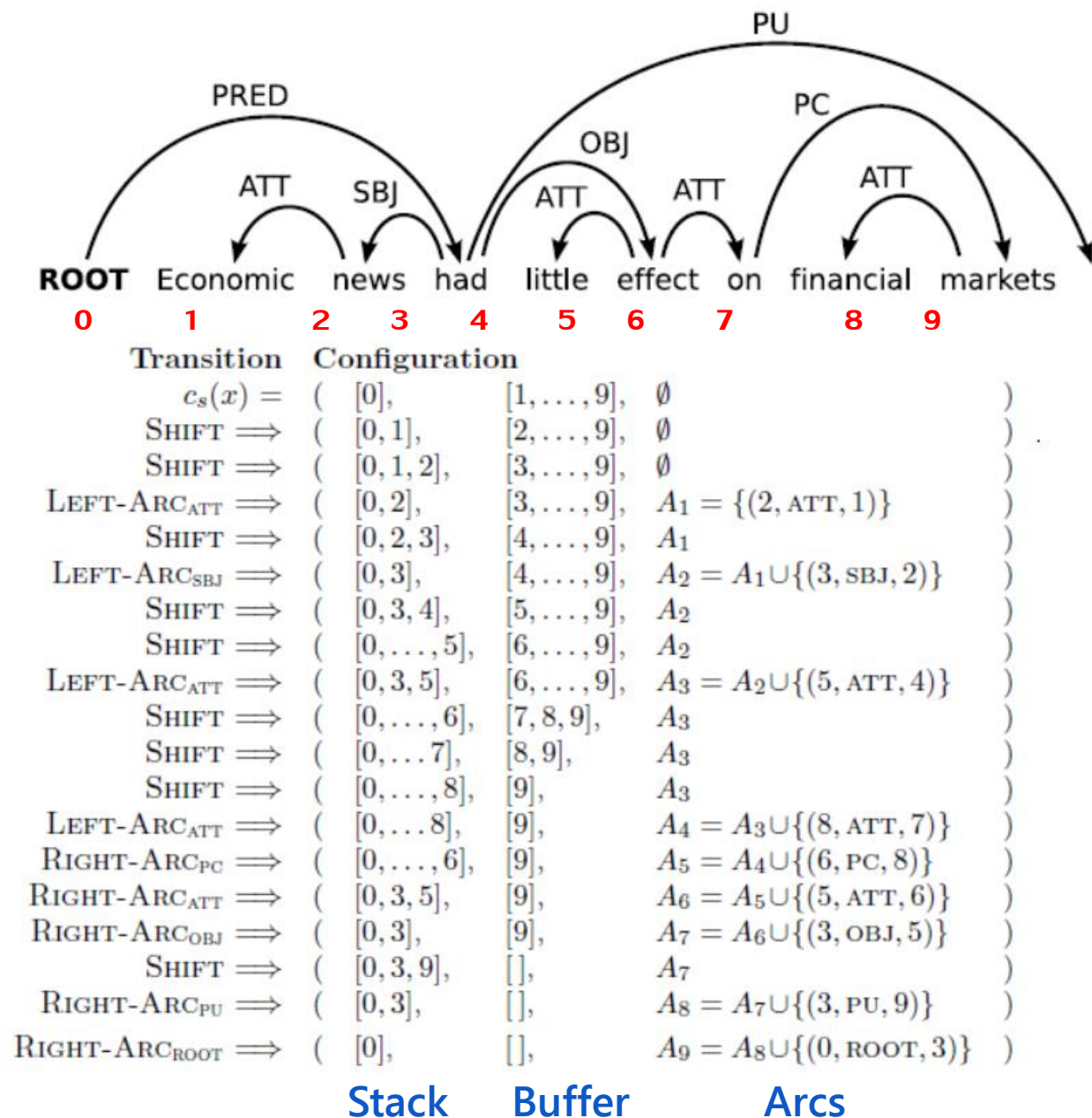
$(\sigma, [i \beta], A) \Rightarrow ([\sigma i], \beta, A)$	(SHIFT)
$([\sigma i j], B, A) \Rightarrow ([\sigma j], B, A \cup \{(j, l, i)\})^1$	(LEFT-ARC _l)
$([\sigma i j], B, A) \Rightarrow ([\sigma i], B, A \cup \{(i, l, j)\})$	(RIGHT-ARC _l)

¹ Permitted only if $i \neq 0$.

The notation $[\sigma|i]$ (for the stack) denotes a right-headed list with head i and tail σ .

The notation $[j|\beta]$ (for the buffer) denotes a left-headed list with head j and tail β .

Transition Sequence For Arc-Standard System



How Arc-Standard Transition works?

$[\text{ROOT}]_{\Sigma}$ [Economic, news, had, little, effect, on, financial, markets, .] $_{\mathcal{B}}$

ROOT Economic news had little effect on financial markets .

How Arc-Standard Transition works?

[ROOT, Economic] Σ [news, had, little, effect, on, financial, markets, .] B

ROOT Economic news had little effect on financial markets .

How Arc-Standard Transition works?


[ROOT, Economic, news] Σ [had, little, effect, on, financial, markets, .] B

ROOT Economic news had little effect on financial markets .

How Arc-Standard Transition works?

[ROOT, news]_Σ [had, little, effect, on, financial, markets, .]_B

ROOT Economic news had little effect on financial markets .



The diagram illustrates an Arc-Standard Transition (ATT) between the words 'Economic' and 'news' in the sentence 'Economic news had little effect on financial markets .'. A curved arrow labeled 'ATT' points from 'Economic' to 'news', indicating a transition between these two words.

How Arc-Standard Transition works?

[ROOT, news, had]_Σ [little, effect, on, financial, markets, .]_B

ROOT Economic news had little effect on financial markets .



The diagram illustrates an Arc-Standard Transition (ATT) between the words 'Economic' and 'news' in the sentence 'Economic news had little effect on financial markets .'. A curved arrow labeled 'ATT' points from 'Economic' to 'news', indicating a transition between these two words.

How Arc-Standard Transition works?

[ROOT, had]_Σ [little, effect, on, financial, markets, .]_B



How Arc-Standard Transition works?

[ROOT, had, little]_Σ [effect, on, financial, markets, .]_B



How Arc-Standard Transition works?

[ROOT, had, little, effect]_Σ [on, financial, markets, .]_B



How Arc-Standard Transition works?

[ROOT, had, effect]_Σ [on, financial, markets, .]_B



How Arc-Standard Transition works?

[ROOT, had, effect, on]_Σ [financial, markets, .]_B



How Arc-Standard Transition works?

[ROOT, had, effect, on, financial]_Σ [markets, .]_B



How Arc-Standard Transition works?

[ROOT, had, effect, on, financial, markets]_Σ [.]_B



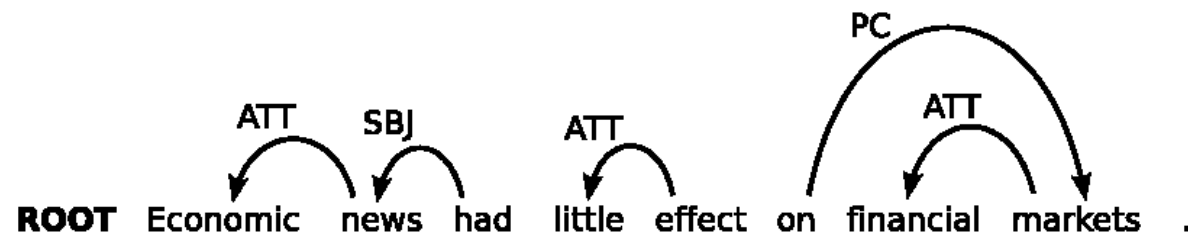
How Arc-Standard Transition works?

[ROOT, had, effect, on, markets]_Σ [.]_B



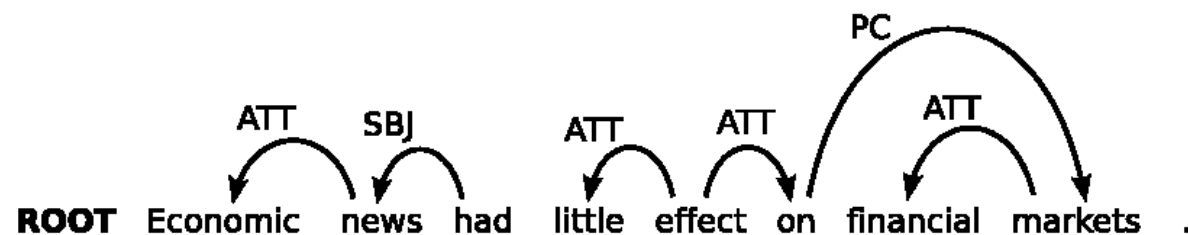
How Arc-Standard Transition works?

[ROOT, had, effect, on]_Σ [.]_B



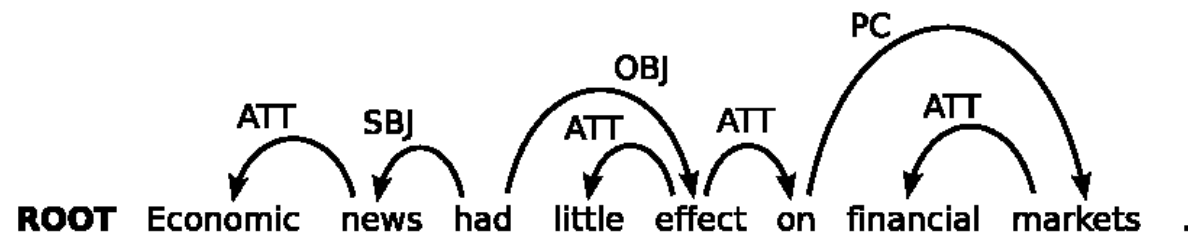
How Arc-Standard Transition works?

[ROOT, had, effect]_Σ [.]_B



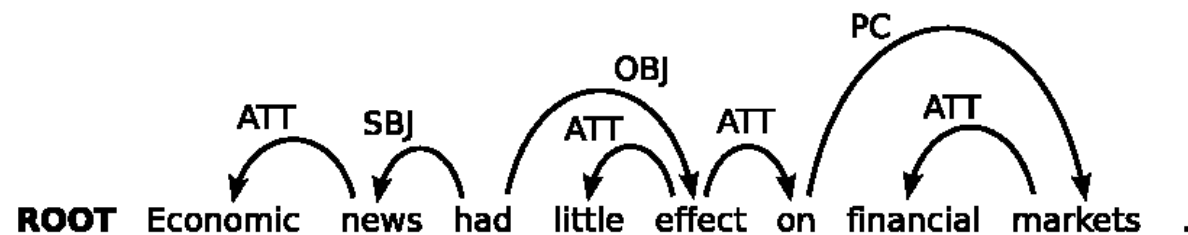
How Arc-Standard Transition works?

[ROOT, had]_Σ [.]_B

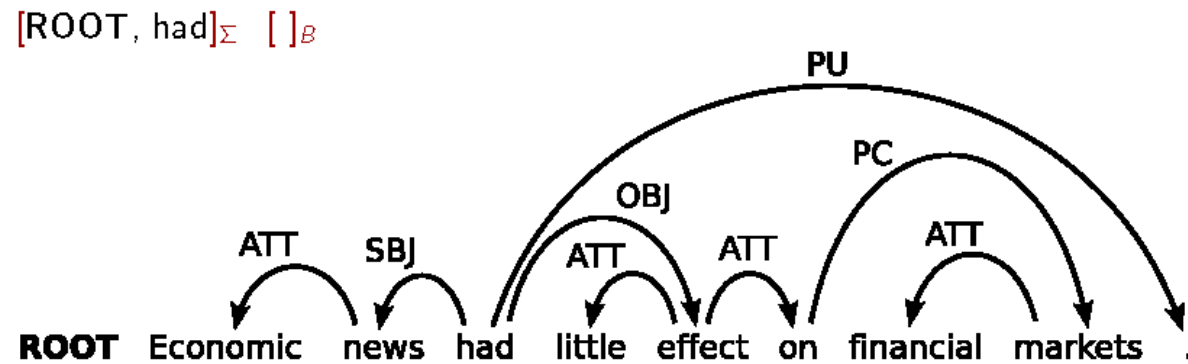


How Arc-Standard Transition works?

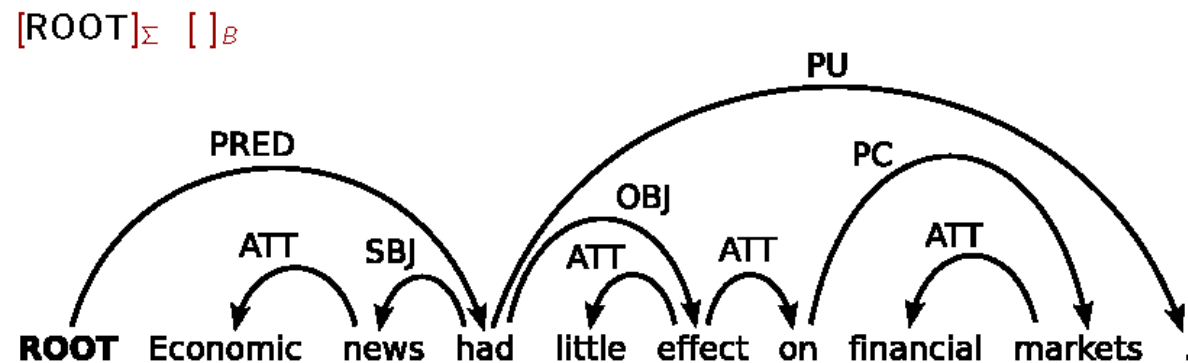
[ROOT, had, .]_Σ []_B



How Arc-Standard Transition works?

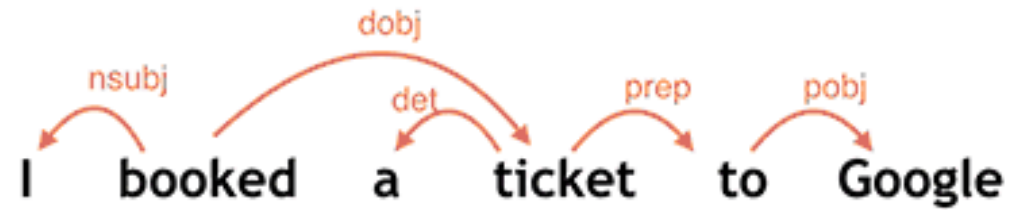


How Arc-Standard Transition works?



Dependency Parsing?

Dependency Parsing





Question & Answer

References

- A Fast and Accurate Dependency Parser using Neural Networks
 - By Danqi Chen, Christopher D. Manning (Stanford)
- Structured Training for Neural Network Transition-Based Parsing
 - By David Wiess, Chris Alberti, Michael Collins, Slav Petrov (Google)
- Globally Normalized Transition-Based Neural Networks
 - By Daniel Andor, Chris Alberti, David Wiess, ... (Google)
- Transition-Based Parsing
 - class material by Joakim Nivre



Thank you