

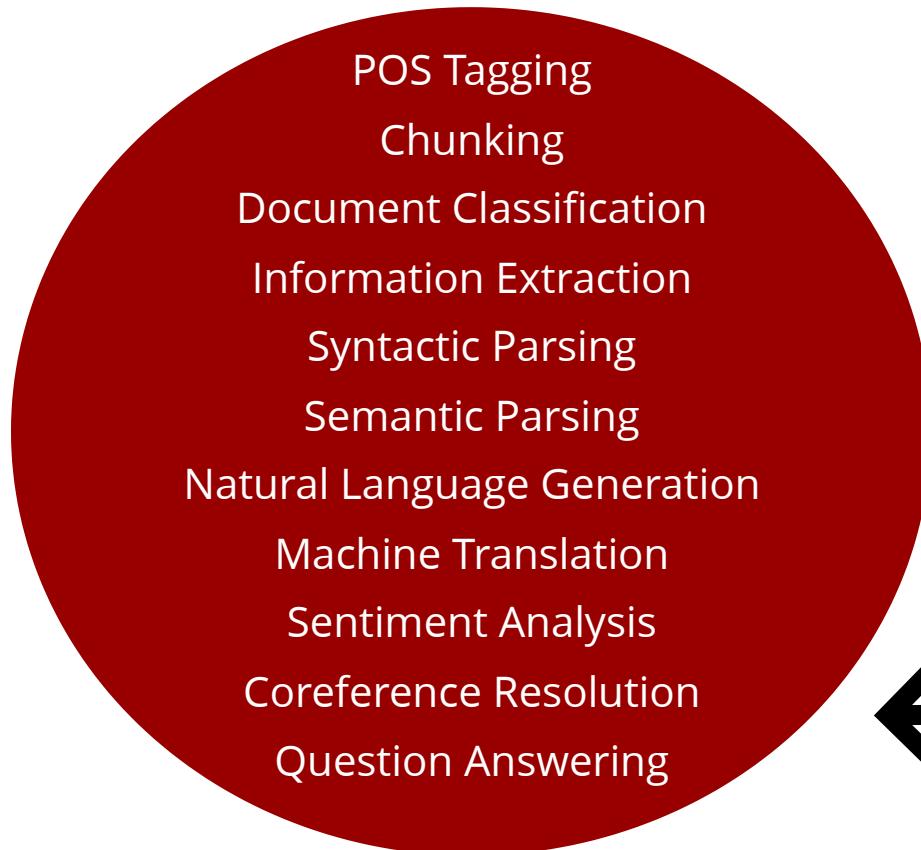
50.040

Natural Language Processing

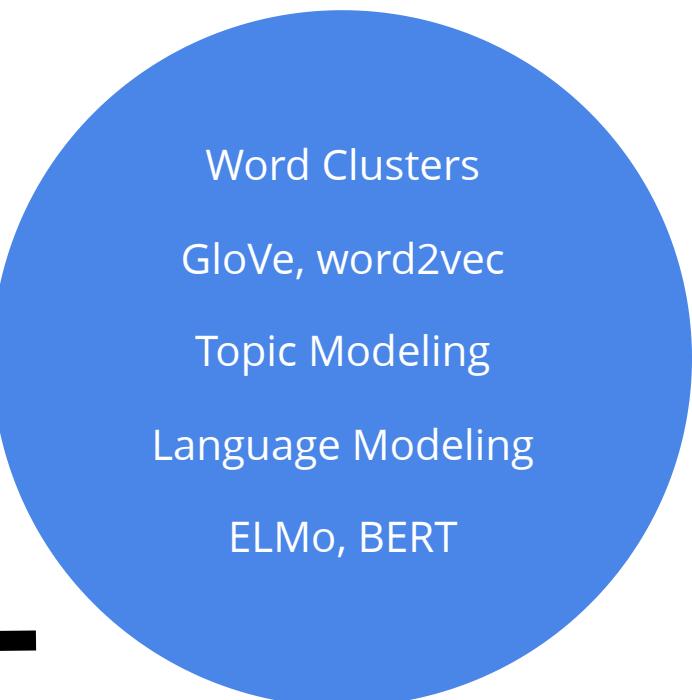
Lu, Wei



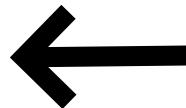
Tasks in NLP



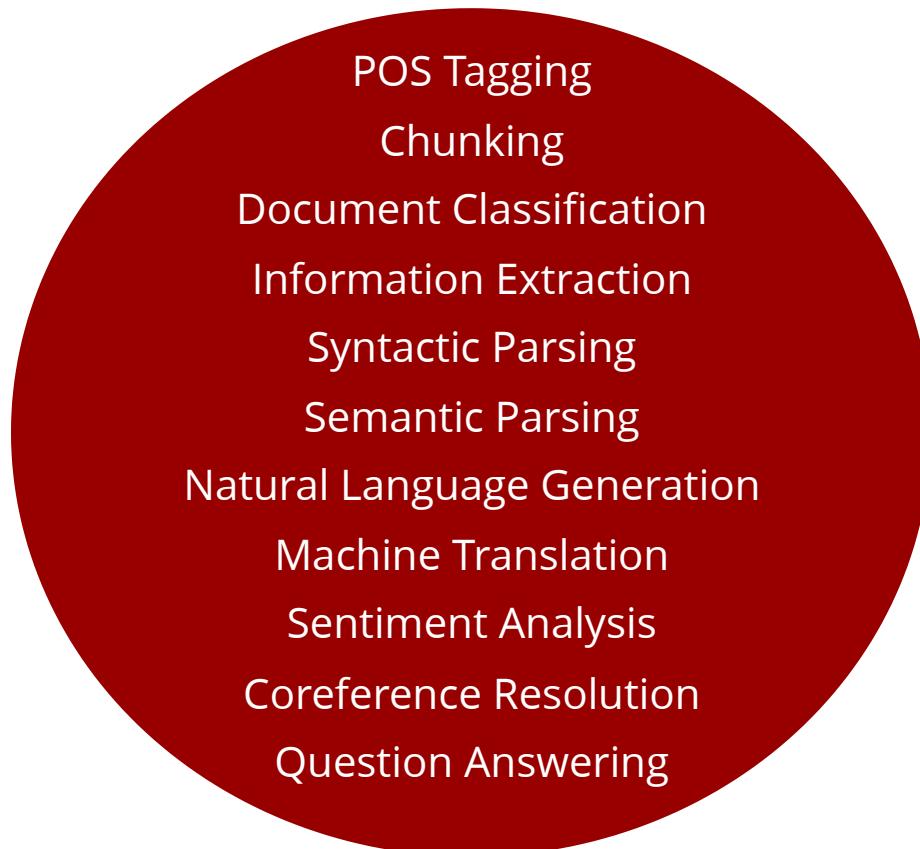
Supervised



Unsupervised

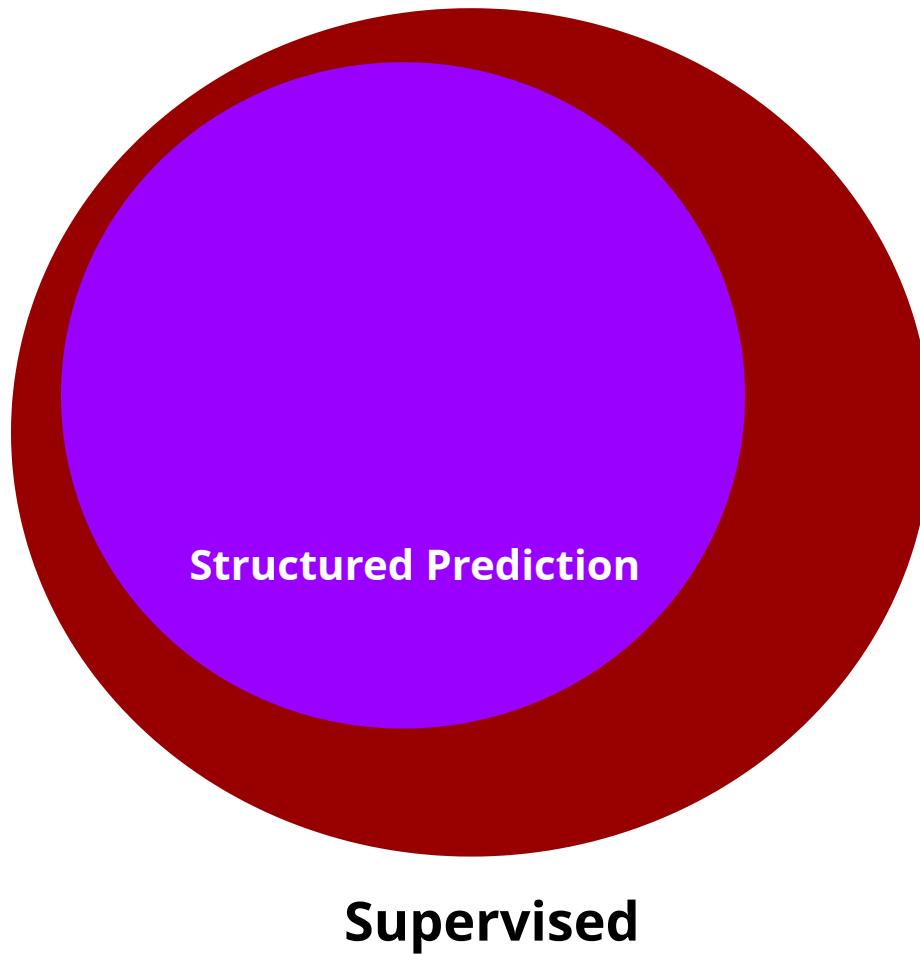


Tasks in NLP



Supervised

Tasks in NLP



Structured Prediction

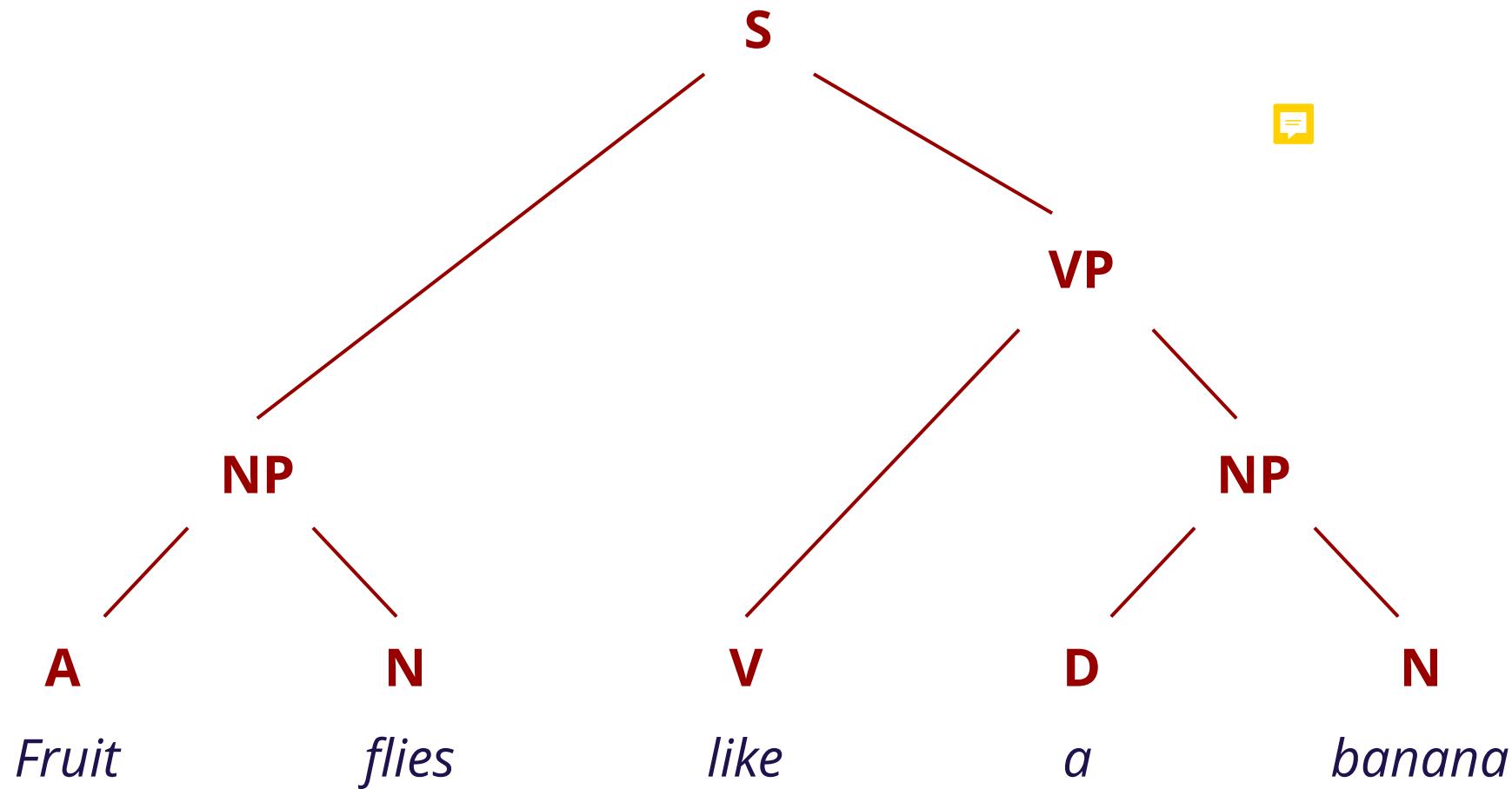
Part-of-Speech Tagging



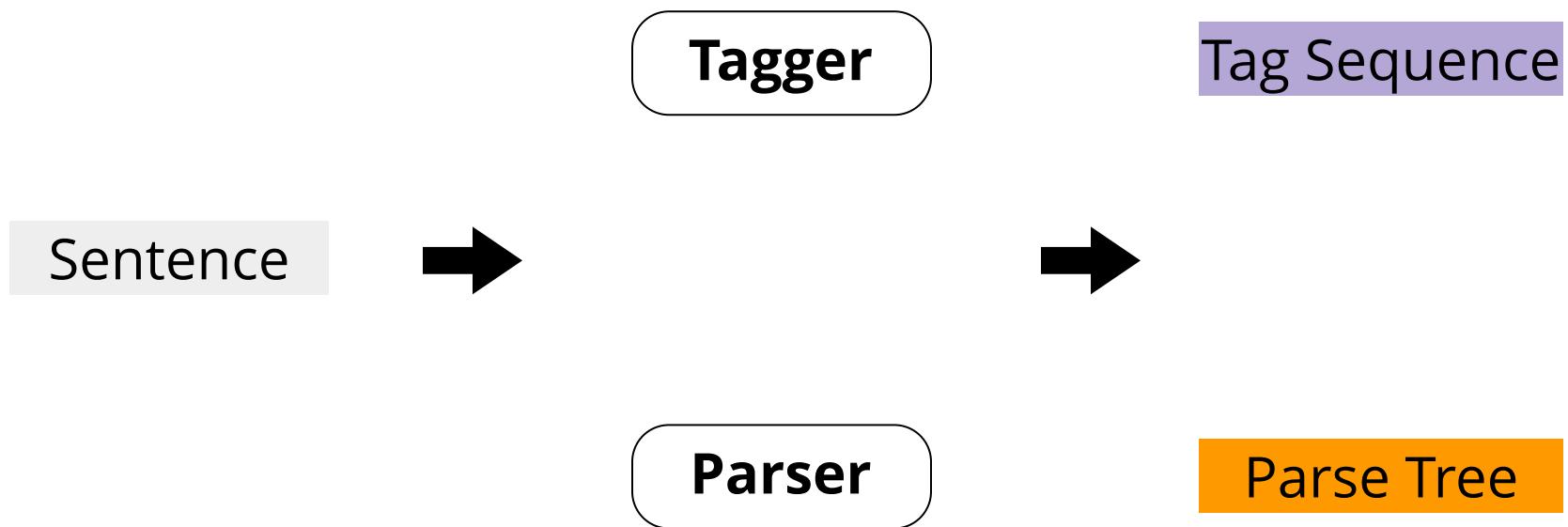
A	N	V	D	N
<i>Fruit</i>	<i>flies</i>	<i>like</i>	<i>a</i>	<i>banana</i>

Structured Prediction

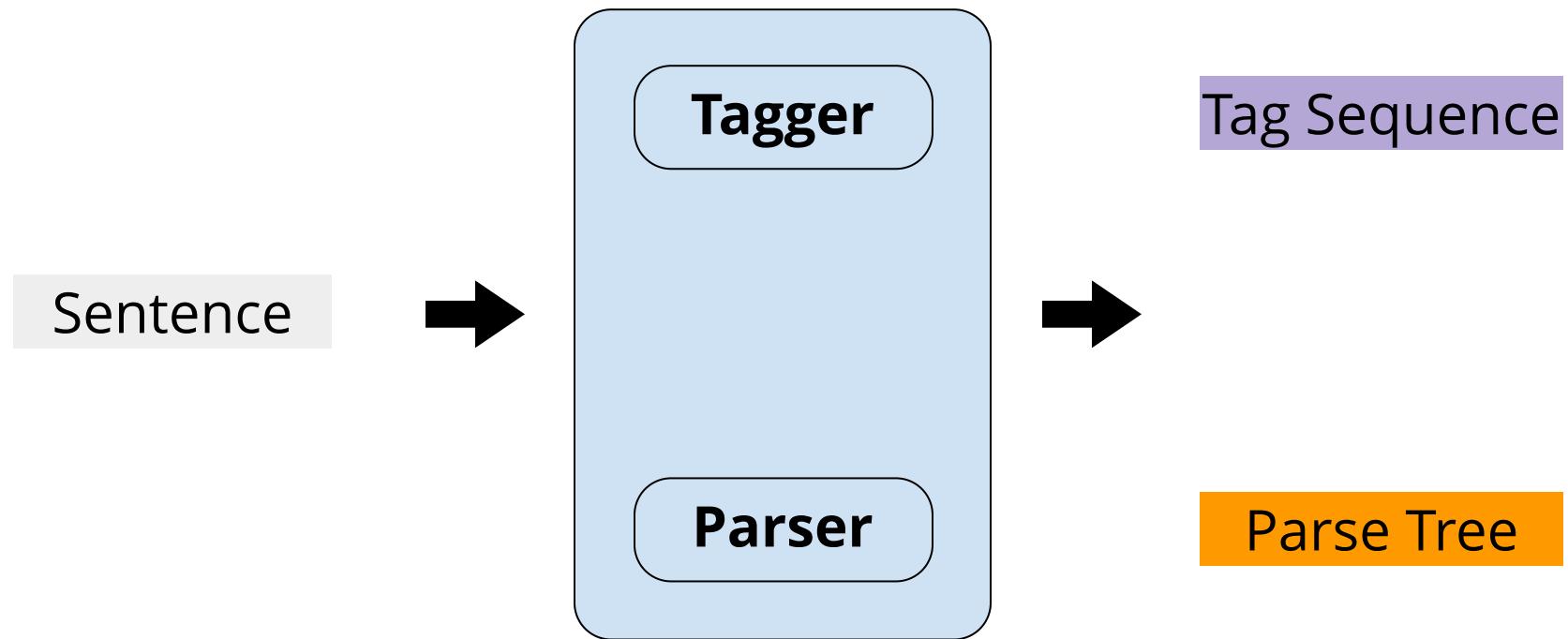
Constituency Parsing



Structured Prediction



Structured Prediction



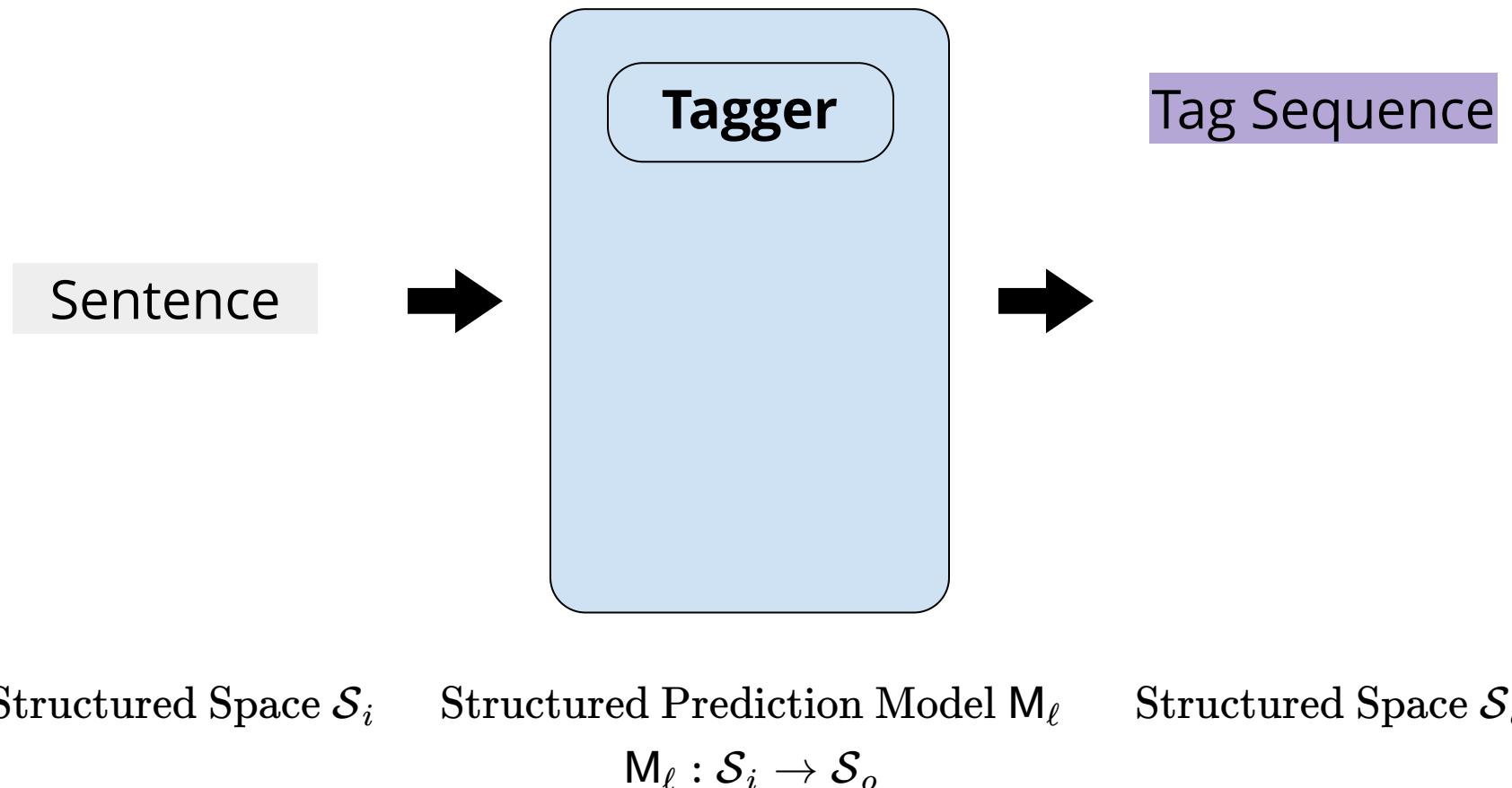
Structured Space \mathcal{S}_i

Structured Prediction Model M

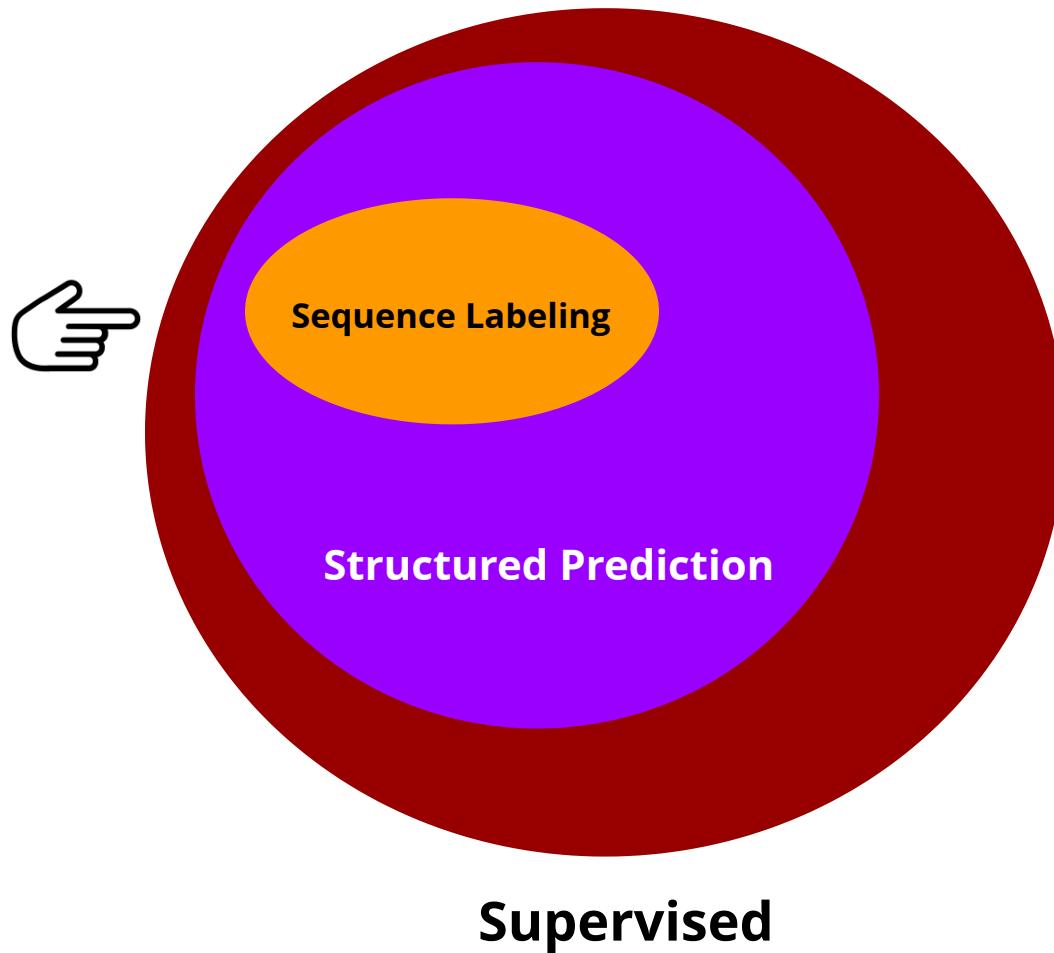
$$M : \mathcal{S}_i \rightarrow \mathcal{S}_o$$

Structured Space \mathcal{S}_o

Sequence Labeling



Tasks in NLP



Sequence Labeling

Definition

A type of machine learning task
that involves assigning one
discrete label to each member of
a sequence of input tokens.

Sequence Labeling

Definition

A type of machine learning task that involves assigning one **discrete label** to each member of a sequence of input tokens.

Example Tasks

Part-of-Speech Tagging

Noun Phrase Chunking

Named Entity Recognition

Chinese Word Segmentation

Part-of-Speech Tagging

Assigning each word token a part-of-speech (POS) tag, where a POS refers to a category of words which have similar grammatical properties.

A	N	V	D	N
<i>Fruit</i>	<i>flies</i>	<i>like</i>	<i>a</i>	<i>banana</i>

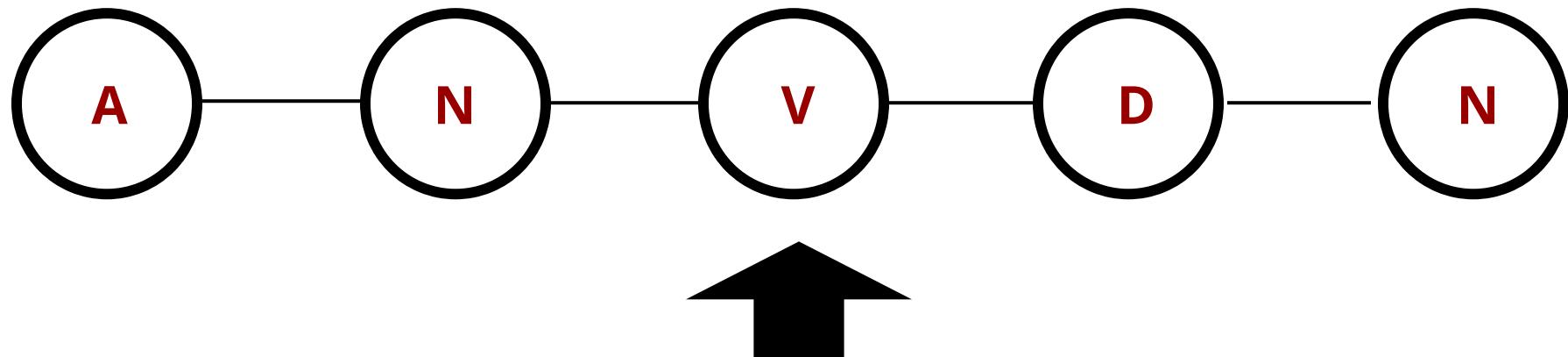
Part-of-Speech Tagging

For English, there are 9 main part of speech (POS): noun, verb, article, adjective, preposition, pronoun, adverb, conjunction, and interjection.

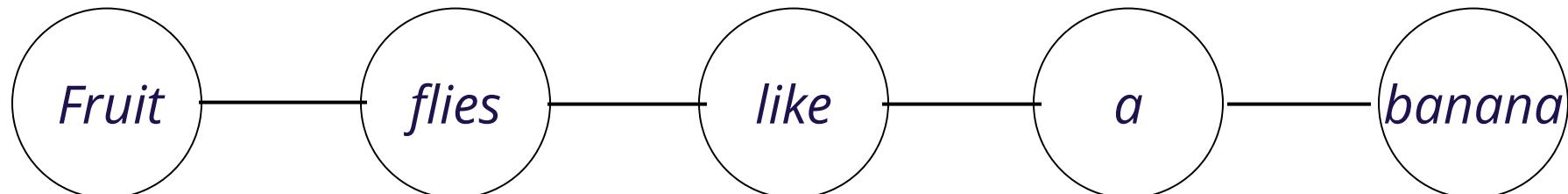
However, in practice there are far more fine-grained POS tags (50 - 150).

A	N	V	D	N
<i>Fruit</i>	<i>flies</i>	<i>like</i>	<i>a</i>	<i>banana</i>

Part-of-Speech Tagging



Tagger \mathcal{M}_ℓ



Brill Tagger (Brill, 1993)

Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part-of-Speech Tagging

Eric Brill*

The Johns Hopkins Univ

Recently, there has been a real encoding of linguistic knowledge as a method of providing a though corpus-based approach to processing, it is often the case that modelling indirectly in large understand and improve the In this paper, we will describe knowledge. This approach has and more direct fashion with of this learning method appl

Transformation-Based Error-Driven Learning

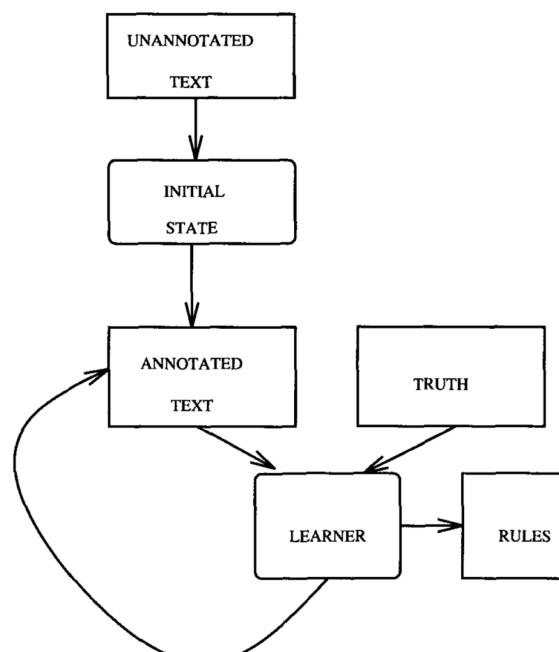
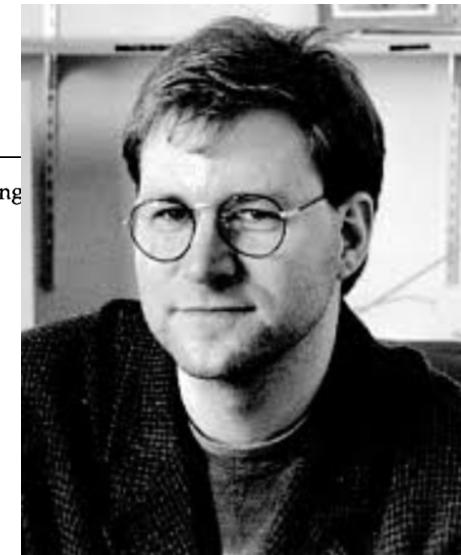


Figure 1
Transformation-Based Error-Driven Learning.

Transformation-Based Error-Driven Learning

Rules are of the form:

Change tag from *a* to *b* when *condition*



```
1. apply initial-state annotator to corpus
2. while transformations can still be found do
3.   for from_tag = tag1 to tagn
4.     for to_tag = tag1 to tagn
5.       for corpus_position = 1 to corpus_size
6.         if (correct_tag(corpus_position) == to_tag)
7.           && current.tag(corpus_position) == from_tag)
8.             num_good_transformations(tag(corpus_position - 1))++
9.         else if (correct_tag(corpus_position) == from_tag)
10.           && current.tag(corpus_position) == from_tag)
11.             num_bad_transformations(tag(corpus_position - 1))++
12.         find maxT (num_good_transformations(T) - num_bad_transformations(T))
13.         if this is the best-scoring rule found yet then store as best rule:
14.           Change tag from from_tag to to_tag if previous tag is T
15.         apply best rule to training corpus
16.         append best rule to ordered list of transformations
```

Figure 3
Pseudocode for learning transformations.

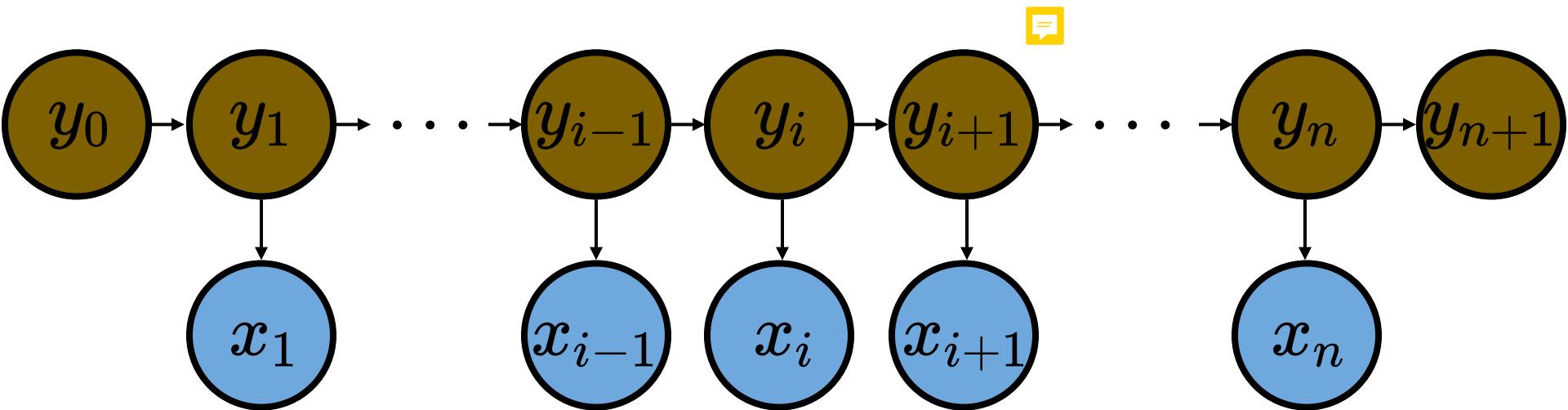
#	From	To	Condition
1	NN	VB	Previous tag is TO
2	VBP	VB	One of the previous three tags is MD
3	NN	VB	One of the previous two tags is MD
4	VB	NN	One of the previous two tags is DT
5	VBD	VBN	One of the previous three tags is VBZ
6	VBN	VBD	Previous tag is PRP
7	VBN	VBD	Previous tag is NNP
8	VBD	VBN	Previous tag is VBD
9	VBP	VB	Previous tag is TO
10	POS	VBZ	Previous tag is PRP
11	VB	VBP	Previous tag is NNS
12	VBD	VBN	One of previous three tags is VBP
13	IN	WDT	One of next two tags is VB
14	VBD	VBN	One of previous two tags is VB
15	VB	VBP	Previous tag is PRP
16	IN	WDT	Next tag is VBZ
17	IN	DT	Next tag is NN
18	JJ	NNP	Next tag is NNP
19	IN	WDT	Next tag is VBD
20	JJR	RBR	Next tag is JJ

Figure 4
The first 20 nonlexicalized transformations.

Hidden Markov Model

Generative Approach

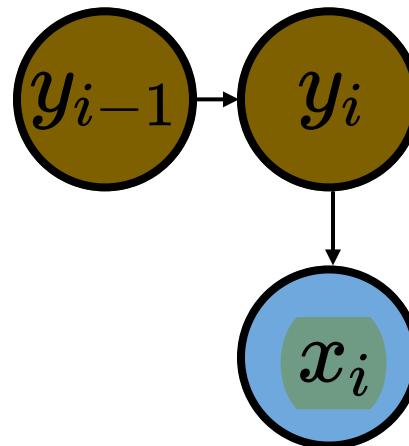
$$\begin{aligned} p(x_1, \dots, x_n, y_0, y_1, \dots, y_n, y_{n+1}) \\ = \prod_{i=1}^{n+1} p(y_i | y_{i-1}) \prod_{i=1}^n p(x_i | y_i) \end{aligned}$$



Hidden Markov Model

One Limitation

$$p(y_i|y_{i-1})p(x_i|y_i)$$

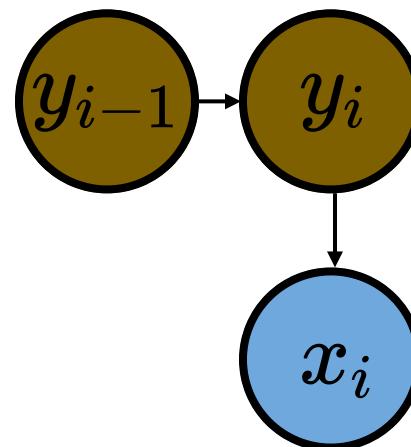


Does a word that ends with "-ed" or "-ly" likely to reveal its POS information?

Hidden Markov Model

One Limitation

$$p(y_i|y_{i-1})p(x_i|y_i)$$



How to capture the rich linguistic features within the inputs?

Conditional Model

Recall what we did in
classification?

We can capture features
with a discriminative
approach!

$$p(y_1, \dots, y_n \mid x_1, \dots, x_n)$$

Conditional Model

$$= \prod_{i=1}^n p(y_i \mid x_1, \dots, x_i, \dots, x_n, y_1, \dots, y_{i-1})$$

Conditional Model

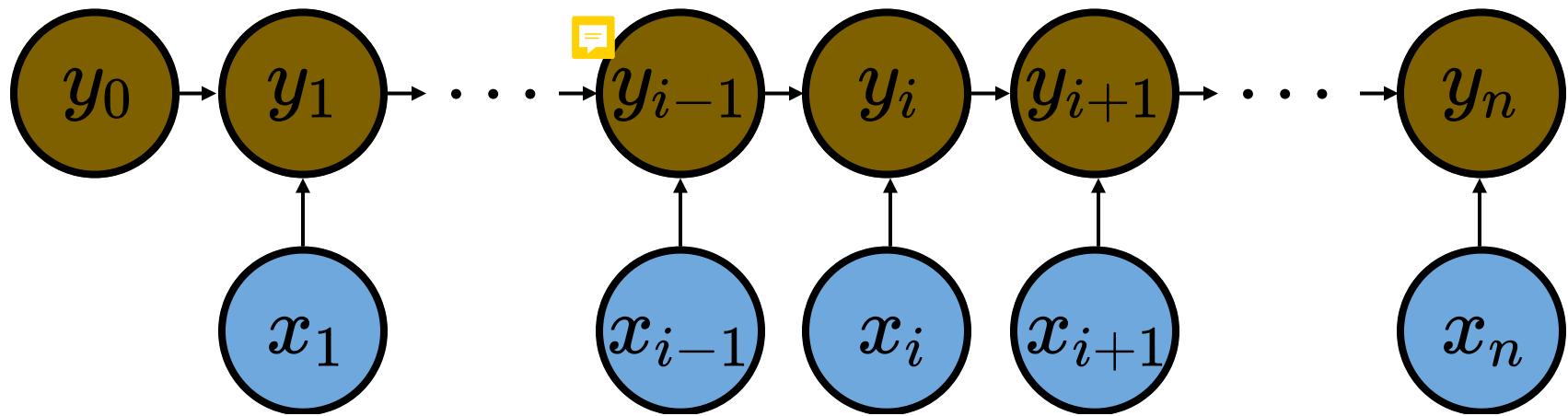
$$\begin{aligned} & p(y_1, \dots, y_n \mid x_1, \dots, x_n) \\ = & \prod_{i=1}^n p(y_i \mid x_1, \dots, x_i, \dots, x_n, y_1, \dots, y_{i-1}) \\ = & \prod_{i=1}^n p(y_i \mid \cancel{x_1}, \cancel{\dots}, x_i, \cancel{\dots}, \cancel{x_n}, \cancel{y_1}, \cancel{\dots}, y_{i-1}) \\ = & \prod_{i=1}^n p(y_i \mid x_i, y_{i-1}) \end{aligned}$$



Conditional Model

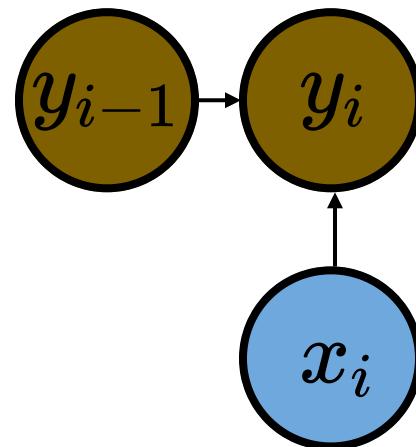
$$\begin{aligned} & p(y_1, \dots, y_n \mid x_1, \dots, x_n) \\ = & \prod_{i=1}^n p(y_i \mid x_1, \dots, x_i, \dots, x_n, y_1, \dots, y_{i-1}) \\ = & \prod_{i=1}^n p(y_i \mid \cancel{x_1}, \cancel{\dots}, \cancel{x_i}, \cancel{\dots}, \cancel{x_n}, \cancel{y_1}, \cancel{\dots}, \cancel{y_{i-1}}) \\ = & \prod_{i=1}^n p(y_i \mid x_i, y_{i-1}) \end{aligned}$$

↑
Markov Assumptions



Conditional Model

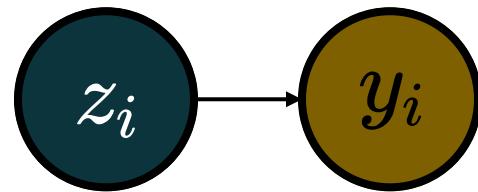
$$p(y_i \mid x_i, y_{i-1})$$



How to parameterize this conditional probability?

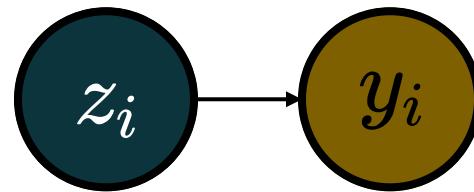
Conditional Model

$p(y_i \mid x_i, y_{i-1}) = p(y_i \mid z_i)$, where $z_i = \langle x_i, y_{i-1} \rangle$



Conditional Model

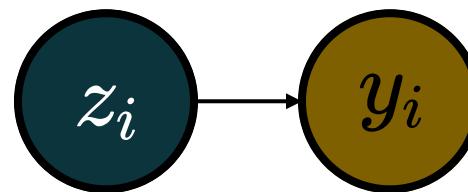
$p(y_i \mid x_i, y_{i-1}) = p(y_i \mid z_i)$, where $z_i = \langle x_i, y_{i-1} \rangle$



What if the discrete variable y_i is binary?

Logistic Regression

$p(y_i \mid x_i, y_{i-1}) = p(y_i \mid z_i)$, where $z_i = \langle x_i, y_{i-1} \rangle$



$$p(y_i \mid z_i) \propto \exp(-(\mathbf{f}(z_i) \cdot \boldsymbol{\theta} + \theta_0)y_i)$$

↑ ↑
feature vector binary variable

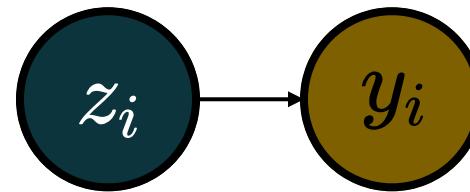
weights/parameters



What if the discrete variable y_i is no longer binary?

Softmax Regression

$p(y_i \mid x_i, y_{i-1}) = p(y_i \mid z_i)$, where $z_i = \langle x_i, y_{i-1} \rangle$



$$p(y_i \mid z_i) \propto \exp(f(z_i, y_i) \cdot \theta)$$

weights/parameters

feature vector

↗

↗

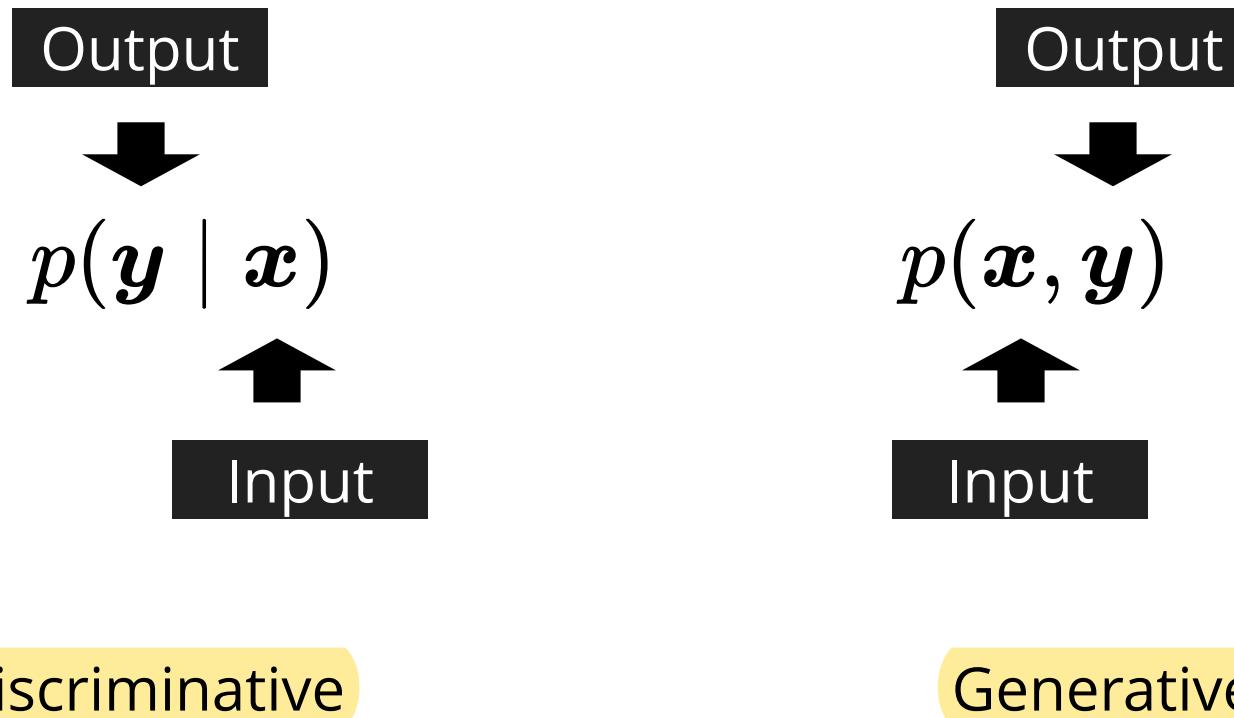


Why can we parameterize the probability this way?



Two Types of Models

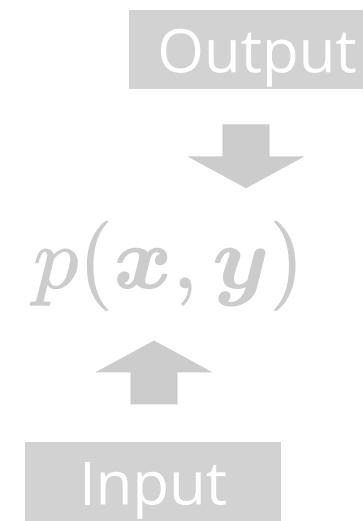
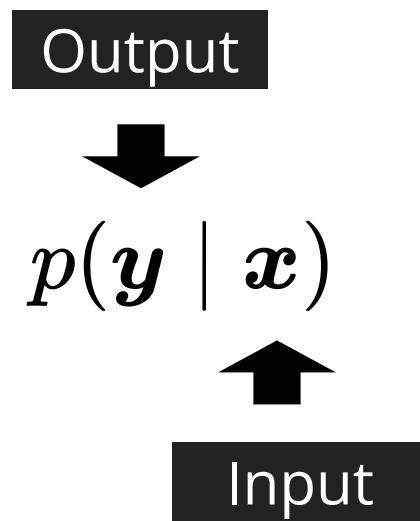
From now, let us use x and y for input and output respectively.



(no need to explicitly model input) (need to explicitly model input)

Discriminative Model

From now, let us use x and y for input and output respectively.



Discriminative

(no need to explicitly model input)

Generative

(need to explicitly model input)

Maximum Entropy

A Maximum Entropy Approach to Natural Language Processing

Adam L. Berger[†]
Columbia University

Vincent J. Della Pietra[‡]
Renaissance Technologies

Stephen A. Della Pietra[‡]
Renaissance Technologies

The concept of maximum entropy can be traced back along multiple threads to Biblical times. Only recently, however, have computers become powerful enough to permit the widespread application of this concept to real world problems in statistical estimation and pattern recognition. In this paper, we describe a method for statistical modeling based on maximum entropy. We also describe a maximum-likelihood approach for automatically constructing maximum entropy models and describe how to implement this approach efficiently, using as examples several problems in natural language processing.



Adam L. Berger, Vincent J. Della Pietra, and Stephen A. Della Pietra.

A Maximum Entropy Approach to Natural Language Processing. Computational Linguistics 22, 1 (March 1996), 39-732

Entropy

Degree of Uncertainty/Surprise

Entropy of a probability distribution $p(v)$
is defined as follows:

$$\mathcal{H}(p(v)) = - \sum_v p(v) \log p(v)$$

It measures the degree of uncertainty associated with the model.



Entropy

Degree of Uncertainty/Surprise

Entropy of a probability distribution $p(v)$
is defined as follows:

$$\mathcal{H}(p(v)) = - \sum_v p(v) \log p(v)$$

It measures the degree of uncertainty associated with the model.

When the number of discrete outputs is fixed, the uniform distribution gives the highest entropy!

Maximum Entropy

The entropy of the conditional model of our interest is defined as:

$$\mathcal{H}(p(\mathbf{y}|\mathbf{x})) = - \sum_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}) \log p(\mathbf{y}|\mathbf{x})$$

Our goal is to find the distribution that maximizes the entropy:

$$p^*(\mathbf{y}|\mathbf{x}) = \max_{p(\mathbf{y}|\mathbf{x})} \mathcal{H}(p(\mathbf{y}|\mathbf{x}))$$



What are the constraints for this optimization problem?



Maximum Entropy

$$p^*(\mathbf{y}|\mathbf{x}) = \max_{p(\mathbf{y}|\mathbf{x})} \mathcal{H}(p(\mathbf{y}|\mathbf{x}))$$

$$\max_{p(\mathbf{y}|\mathbf{x})} \mathcal{H}(p(\mathbf{y}|\mathbf{x}))$$

subject to:

1. It has to be a valid probability distribution:

$$\sum_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}) = 1$$

2. While we prefer a flat/uniform distribution, we shall respect the empirical features that we observed:

$$\mathbf{E}_{p(\mathbf{y}|\mathbf{x})}[f_k(\mathbf{x}, \mathbf{y})] = \mathbf{E}_{\hat{p}(\mathbf{y}|\mathbf{x})}[f_k(\mathbf{x}, \mathbf{y})]$$



The distribution we are
looking for



The empirical
distribution

Maximum Entropy



$$\begin{aligned}\Lambda(p, \lambda) \\ = -\mathcal{H}(p(\mathbf{y}|\mathbf{x})) \\ - \sum_{k=1}^K \lambda_k \left(\mathbf{E}_{p(\mathbf{y}|\mathbf{x})[f_k(\mathbf{x}, \mathbf{y})]} - \mathbf{E}_{\hat{p}(\mathbf{y}|\mathbf{x})[f_k(\mathbf{x}, \mathbf{y})]} \right) \\ - \lambda_0 \left(\sum_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}) - 1 \right)\end{aligned}$$

where λ_k and $\lambda_0 \in (-\infty, +\infty)$

Take the partial derivative and set it to zero:

$$\frac{\partial \Lambda(p, \lambda)}{\partial p(\mathbf{y}|\mathbf{x})} = \log p(\mathbf{y}|\mathbf{x}) + 1 - \sum_{k=1}^K \lambda_k f_k(\mathbf{x}, \mathbf{y}) - \lambda_0 = 0$$

This gives us:

$$p(\mathbf{y}|\mathbf{x}) \propto \exp \left(\sum_{k=1}^K \lambda_k f_k(\mathbf{x}, \mathbf{y}) \right)$$

Maximum Entropy

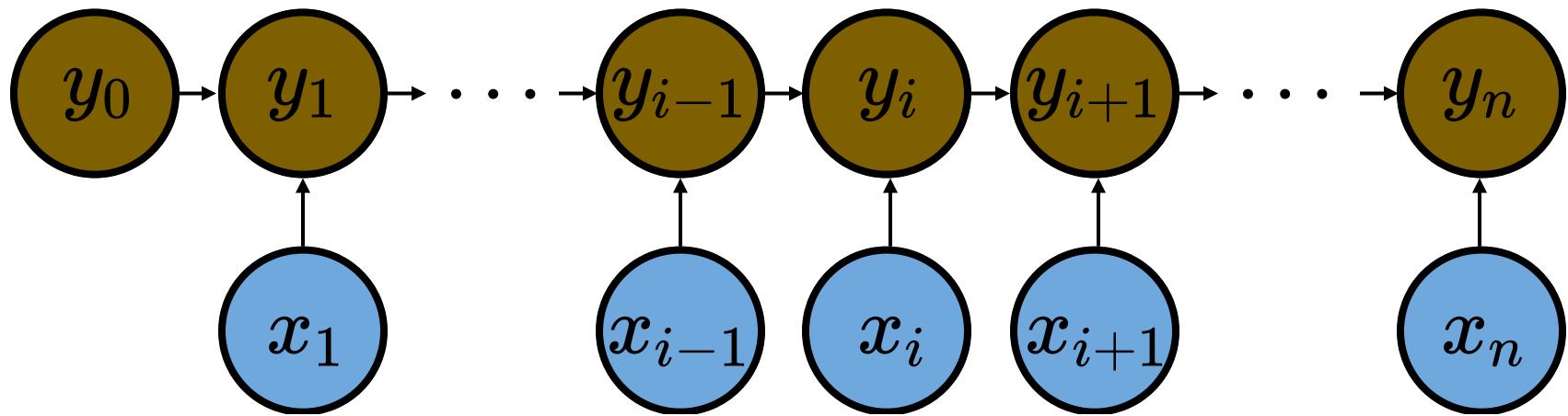
Markov Model

(McCallum et al., 2000)

$$p(y_1, \dots, y_n \mid x_1, \dots, x_n) = \prod_{i=1}^n p(y_i \mid x_i, y_{i-1})$$



$$p(y_i \mid x_i, y_{i-1}) \propto \exp \left(\sum_{k=1}^K \lambda_k f_k(x_i, y_{i-1}, y_i) \right)$$



Part-of-Speech Tagging

Summary

Model	Pros	Cons
Transformation-based Model (Brill's Tagger)	Intuitive, interpretable, fast at testing time	Non probabilistic, training can be slow, accuracy not high enough
Hidden Markov Model (HMM)	Probabilistic, efficient training	Unable to exploit linguistic features
Maximum Entropy Markov Model (MEMM)	Probabilistic, able to exploit features, higher accuracy than HMM	Training is slower than HMM

Noun-Phrase Chunking

NP **NP**
Fruit *flies* *like* *a* *banana*

Question

Is it possible to perform
chunking with tagging?

Chunking as Tagging

Essentially, we are assigning a tag to each individual word in the sentence



Can we just use the below tags for this task? Why?



Part of
NP

Part of
NP

Not Part of
NP

Part of
NP

Part of
NP

NP

NP

Fruit

flies

like

a

banana

Chunking

Tagging Schemes

BIO Tagging Scheme

B: Beginning of NP

I: Inner part of NP

O: Outside of NP

BILOU Tagging Scheme

B: Beginning of NP

I: Inner part of NP

L: Last word of NP

O: Outside of NP

U: Unit word NP

B

I

O

B

I

NP

NP

Fruit

flies

like

a

banana

Chunking

Tagging Schemes

BIO Tagging Scheme

B: Beginning of NP

I: Inner part of NP

O: Outside of NP

BILOU Tagging Scheme

B: Beginning of NP

I: Inner part of NP

L: Last word of NP

O: Outside of NP

U: Unit word NP

B

L

O

B

L

NP

NP

Fruit

flies

like

a

banana

Chunking

Chunking as Shallow Parsing

Captures shallow, non-hierarchical syntactic information.

Identifies elementary constituent parts of a sentence.

As a very basic constituency parser with a flat representation

NP

NP

Fruit

flies

like

a

banana

Named Entity Recognition

One of the most fundamental tasks within NLP, a sub-task within information extraction, that seeks to recognize and classify named entity mentions in text into pre-defined semantic categories (e.g., location, organization, geo-political entity, person, facility ...)

ORG

FAC

Singapore University of Technology and Design is close to the Changi Airport.

Question

How to perform named entity
recognition with tagging?

Named Entity Recognition

Tagging Scheme

B-ORG I-ORG I-ORG I-ORG I-ORG L-ORG O O O O B-FAC L-FAC

ORG

FAC

Singapore University of Technology and Design is close to the Changi Airport.

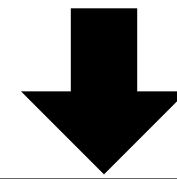
Word Segmentation



近百名后港居民到白沙公园净滩

(早报讯) 后港居民趁星期天早晨展开净滩活动，
为绿化环境尽一份力。

后港基层组织年度净滩活动，号召近百名居民参
加。他们上午8时30分齐聚白沙公园，每人一手拿钳
子，一手拿垃圾袋，在沙滩上捡垃圾。



近 百 名 后 港 居 民 到 白 沙 公 园
净 滩

(早报讯) 后港居民趁星期天早晨展开净滩
活动，为绿化环境尽一份力。

后港基层组织年度净滩活动，号召近百名居民参
加。他们上午8时30分齐聚白沙公园，每人一手拿钳
子，一手拿垃圾袋，在沙滩上捡垃圾。

Word Segmentation

首页 > 即时报道 > 新加坡即时

李总理：来届大选候选人背景多元

发布 / 2019年6月7日 11:43 PM
更新 / 2019年6月8日 5:26 PM
文 / 杨浚鑫

(早报讯) 执政的人民行动党正在为来届大选备战，并已物色许多候选人，身为行动党秘书长的李显龙总理也亲自面试了好些人。

李总理今晚以主宾身份出席第九届“通商中国奖”颁奖晚宴，并在约50分钟的对话会中，谈及我国领导接棒过程和来届大选。他透露，新一批候选人来自不同种族、年龄和背景，其中一些人也曾走过不同的人生道路。

“他们也许在学校表现不佳，到了非政府组织发展事业，并对此深具热情。他们也可能在海外待过一段日子，更了解世界，回国后也更珍惜新加坡。”

行动
(act/action)

党
(party)

行动党
(PAP)

Word Segmentation

首页 > 即时报道 > 新加坡即时

李总理：来届大选候选人背景多元

发布 / 2019年6月7日 11:43 PM
更新 / 2019年6月8日 5:26 PM
文 / 杨浚鑫

f    

中一
(secondary one)
其中
(out of)
一些
(some)
人
(people)

（早报讯）执政的人民行动党正在为来届大选备战，并已物色许多候选人，身为行动党秘书长的李显龙总理也亲自面试了好些人。

李总理今晚以主宾身份出席第九届“通商中国奖”颁奖晚宴，并在约50分钟的对话及我国领导接棒过程和来届大选。他透露，新一批候选人来自不同种族、年岁，其中一些人也曾走过不同的人生道路。

“他们也许在学校表现不佳，到了非政府组织发展事业，并对此深具热情。他们也可能在海外待过一段日子，更了解世界，回国后也更珍惜新加坡。”

Word Segmentation

首页 > 即时报道 > 新加坡即时

李总理：来届大选候选人背景多元

发布 / 2019年6月7日 11:43 PM

更新 / 2019年6月8日 5:26 PM

文 / 杨浚鑫



(早报讯) 执政的人民行动党正在为来届大选备战，并已物色许多候选人，身为行动党秘书长的李显龙总理也亲自面试了好些人。

李总理今晚以主宾身份出席第九届“通商中国奖”颁奖晚宴，并在约50分钟的对话

中，谈及我国领导接棒过程和来届大选。他透露，新一批候选人来自不同种族、年龄和背景，其中一些人也曾走过不同的人生道路。

“他们也许在学校表现不佳，到了非政府组织发展事业，并对此深具热情。他们也可能在海外待过一段日子，更了解世界，回国后也更珍惜新加坡。”

非

(not/African)

政府

(government)

组织

(organization/tissue)

非政府组织

(NGO)

Word Segmentation Tagging Scheme

Consider this Chinese phrase and consider the BILOU scheme:

他到了非政府组织发展事业

Assume we know the gold segmentation is as follows:

他 到了 非 政 府 组 织 发 展 事 业



How shall we tag each Chinese character?

他到了非政府组织发展事业

Word Segmentation Tagging Scheme

Consider this Chinese phrase and consider the BILOU scheme:

他到了非政府组织发展事业

Assume we know the gold segmentation is as follows:

他 到 了 非 政 府 组 织 发 展 事 业

U B L B I I I L B L B L



他到了非政府组织发展事业

Maximum Entropy Markov Model

One Theoretical Limitation

$$p(y_1, \dots, y_n \mid x_1, \dots, x_n) = \prod_{i=1}^n p(y_i \mid x_i, y_{i-1})$$

The model predicts one tag at a time (i.e., performs local normalization).
Since we are modeling sequences, is it possible to predict a complete sequence at a time (i.e., perform global normalization)?

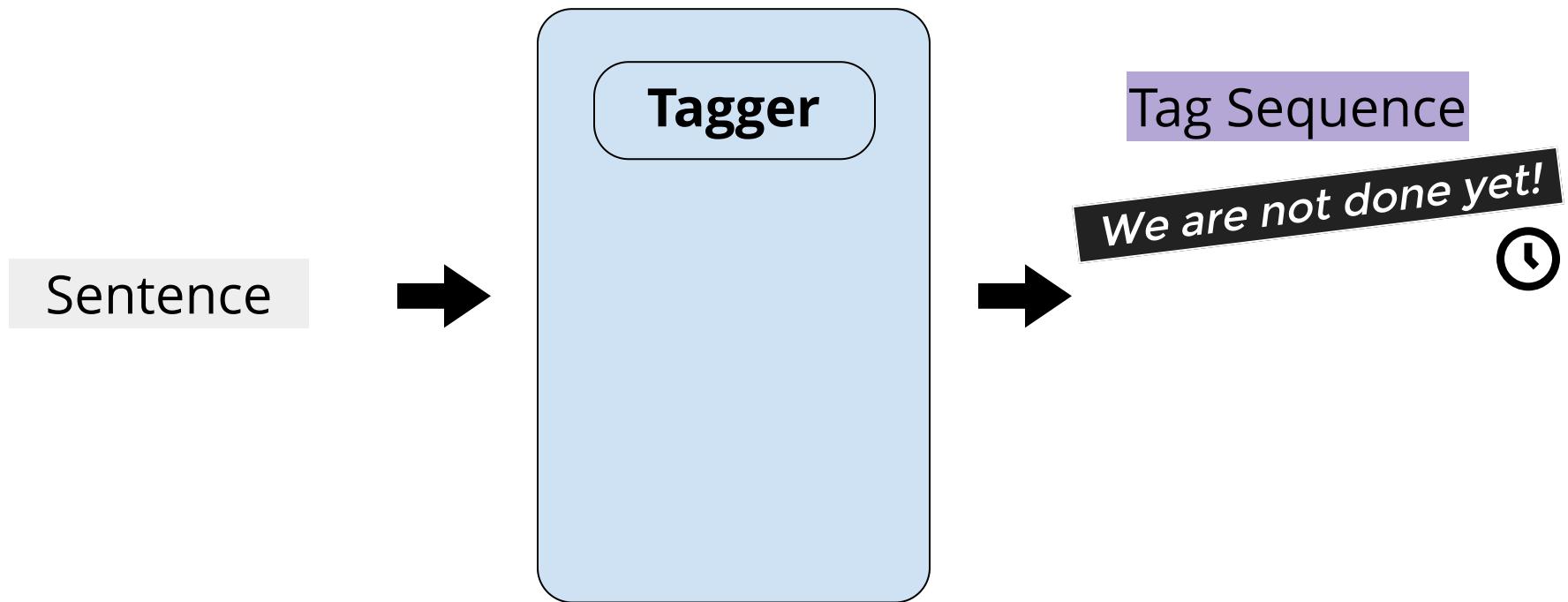


$$p(y_1, \dots, y_n \mid x_1, \dots, x_n) \quad ?$$

$$p(y \mid x) \propto \exp(f(x, y) \cdot \theta)$$

Conditional
Random Fields!

Structured Prediction



Structured Space \mathcal{S}_i

Structured Prediction Model M_ℓ

Structured Space \mathcal{S}_o

$$M_\ell : \mathcal{S}_i \rightarrow \mathcal{S}_o$$