

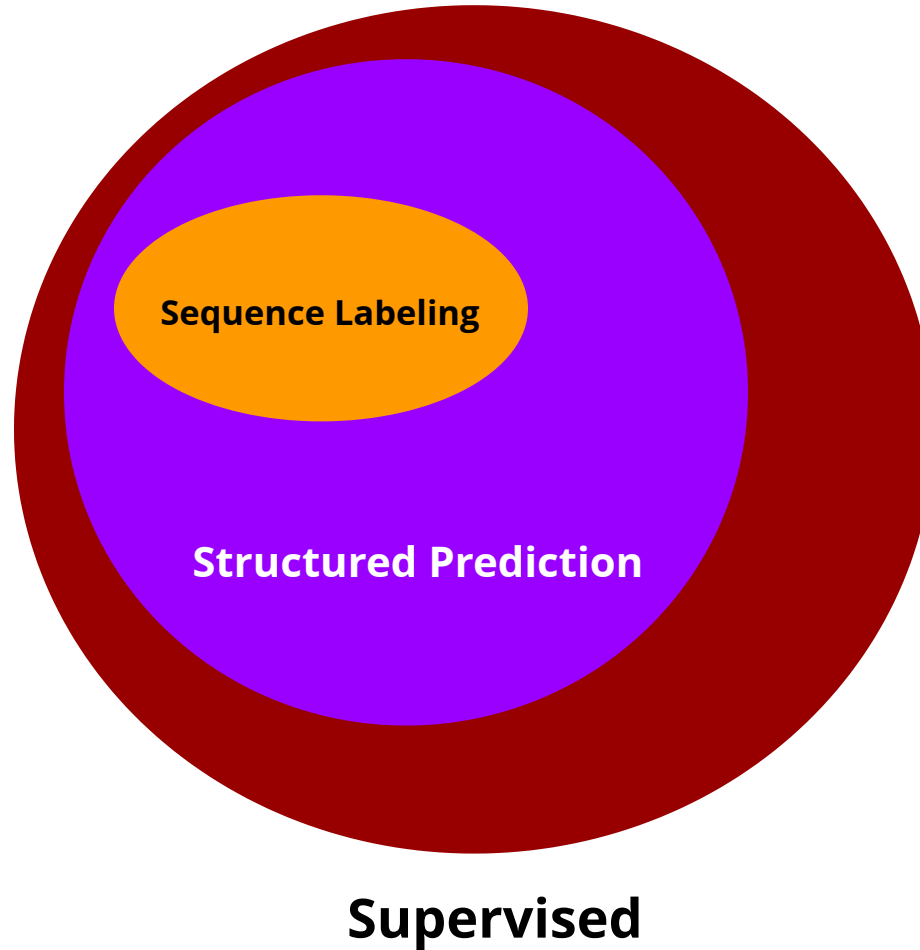
50.040

Natural Language Processing

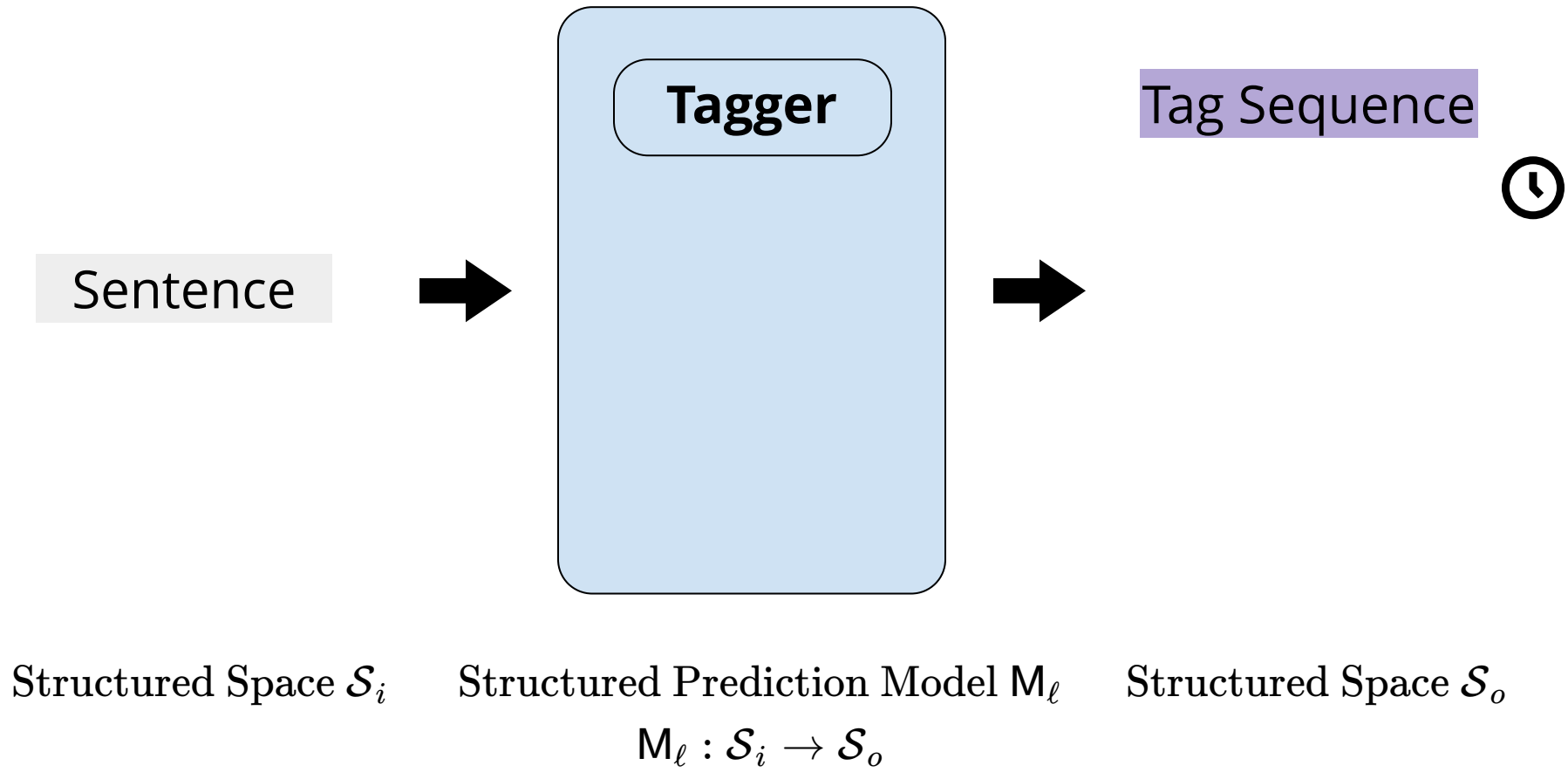
Lu, Wei



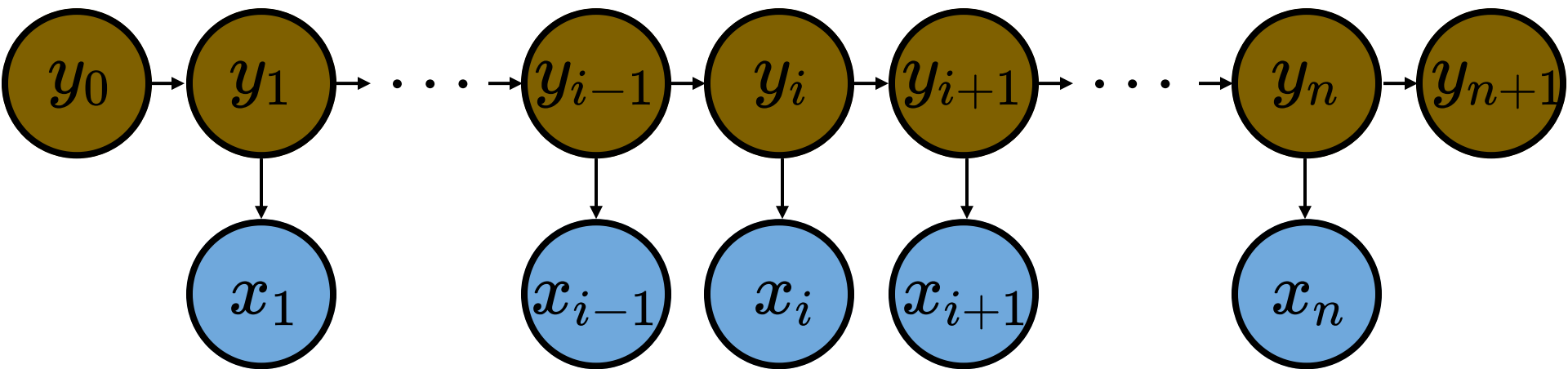
Tasks in NLP



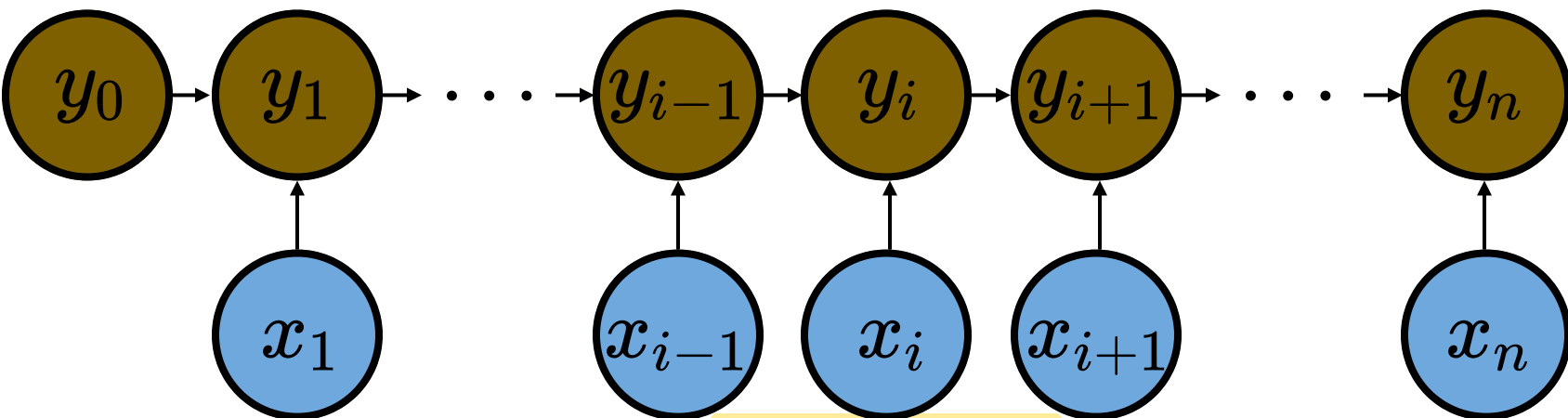
Sequence Labeling



HMM vs MEMM



HMM

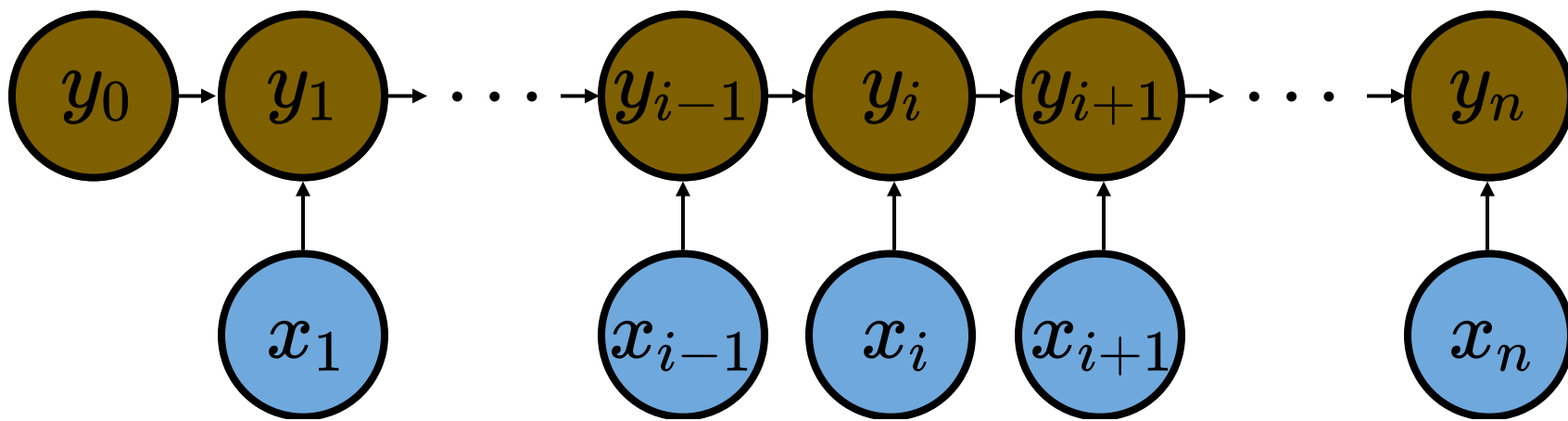


MEMM

Maximum Entropy Markov Model (McCallum et al., 2000)

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

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



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

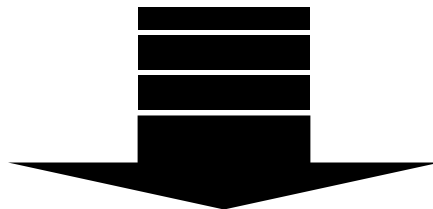


The empirical
distribution

Maximum Entropy

$$p(\mathbf{y} \mid \mathbf{x}_i) = \frac{\exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \boldsymbol{\theta})}{\sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta})}$$

$$\begin{aligned} & \min_{p(\mathbf{y}|\mathbf{x}_i)} \max_{\boldsymbol{\theta}} -\mathcal{H}(p(\mathbf{y}|\mathbf{x}_i)) \\ & - \sum_{k=1}^K \lambda_k \left(\mathbf{E}_{p(\mathbf{y}|\mathbf{x}_i)} [f_k(\mathbf{x}_i, \mathbf{y})] - f_k(\mathbf{x}_i, \mathbf{y}_i) \right) \\ & - \lambda_0 \left(\sum_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}) - 1 \right) \end{aligned}$$



$$\min_{\boldsymbol{\theta}} - \sum_i \log p(\mathbf{y}_i \mid \mathbf{x}_i)$$

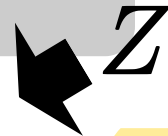
Maximum Entropy

$$\min_{\theta} - \sum_i \log p(\mathbf{y}_i | \mathbf{x}_i)$$

Consider the
 i -th instance



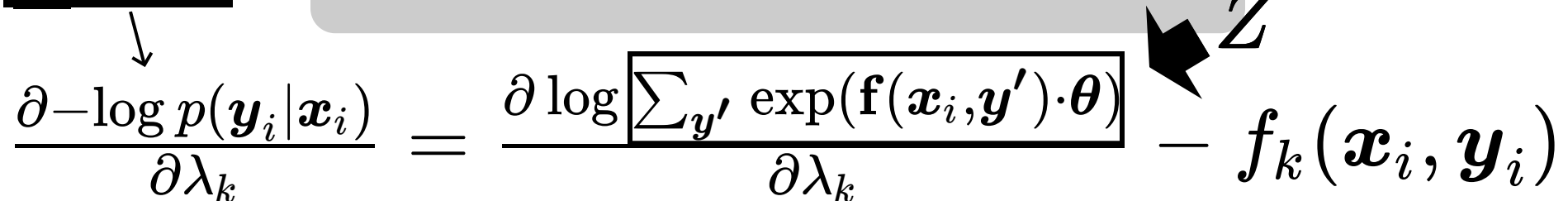
$$\frac{\partial -\log p(\mathbf{y}_i | \mathbf{x}_i)}{\partial \lambda_k} = \frac{\partial \log \left[\sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta}) \right]}{\partial \lambda_k} - f_k(\mathbf{x}_i, \mathbf{y}_i)$$



Maximum Entropy

$$\min_{\theta} - \sum_i \log p(\mathbf{y}_i | \mathbf{x}_i)$$

Consider the
 i -th instance

$$\frac{\partial -\log p(\mathbf{y}_i | \mathbf{x}_i)}{\partial \lambda_k} = \frac{\partial \log \boxed{\sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta})}}{\partial \lambda_k} - f_k(\mathbf{x}_i, \mathbf{y}_i)$$


$$\frac{\partial \log \sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta})}{\partial \lambda_k} = \frac{1}{Z} \cdot \frac{\partial \sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta})}{\partial \lambda_k}$$

$$= \frac{1}{Z} \cdot \sum_{\mathbf{y}'} \frac{\partial \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta})}{\partial \lambda_k}$$

$$= \sum_{\mathbf{y}'} \frac{\exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta})}{Z} f_k(\mathbf{x}_i, \mathbf{y}')$$

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

Maximum Entropy

$$\min_{\theta} - \sum_i \log p(\mathbf{y}_i | \mathbf{x}_i)$$

$$\frac{\partial -\log p(\mathbf{y}_i | \mathbf{x}_i)}{\partial \lambda_k} = \frac{\partial \log \sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta})}{\partial \lambda_k} - f_k(\mathbf{x}_i, \mathbf{y}_i)$$

$$= \dots$$

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



Setting this to zero is exactly one of the two constraints in the optimization problem!

Maximum Entropy

$$\min_{\theta} - \sum_i \log p(\mathbf{y}_i \mid \mathbf{x}_i)$$

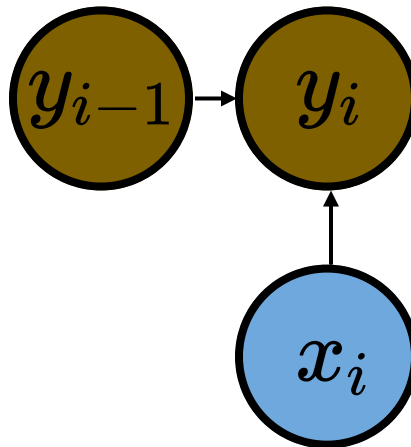
$$\begin{aligned} \frac{\partial -\log p(\mathbf{y}_i \mid \mathbf{x}_i)}{\partial \lambda_k} &= \frac{\partial \log \sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}') \cdot \boldsymbol{\theta})}{\partial \lambda_k} - f_k(\mathbf{x}_i, \mathbf{y}_i) \\ &= \dots \\ &= \mathbf{E}_{p(\mathbf{y}' \mid \mathbf{x}_i)} [f_k(\mathbf{x}_i, \mathbf{y}')] - f_k(\mathbf{x}_i, \mathbf{y}_i) \end{aligned}$$

We can then use (batch) stochastic gradient descent to optimize the objective!

Maximum Entropy Markov Model

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

Each term is a locally trained
Maximum Entropy classifier!





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)?


$$p(y_1, \cdot \mathbf{y} \cdot, y_n \mid x_1, \cdot \mathbf{x} \cdot, x_n)$$
$$p(\mathbf{y} \mid \mathbf{x}) \propto \exp(\mathbf{f}(\mathbf{x}, \mathbf{y}) \cdot \boldsymbol{\theta})$$


Conditional Random Fields

(Lafferty et al., 2001)

Let's now try to model complete sequences

This is now a complete word sequence



$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp(\mathbf{f}(\mathbf{x}, \mathbf{y}) \cdot \boldsymbol{\theta})}{\sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}, \mathbf{y}') \cdot \boldsymbol{\theta})}$$



This is now a complete tag sequence

Can we calculate the denominator efficiently?

Learning in CRF

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp(\mathbf{f}(\mathbf{x}, \mathbf{y}) \cdot \boldsymbol{\theta})}{\sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}, \mathbf{y}') \cdot \boldsymbol{\theta})}$$

$$\min_{\boldsymbol{\theta}} - \sum_i \log p(\mathbf{y}_i \mid \mathbf{x}_i)$$

$$\frac{\partial -\log p(\mathbf{y}_i \mid \mathbf{x}_i)}{\partial \lambda_k} = \underbrace{\mathbf{E}_{p(\mathbf{y}' \mid \mathbf{x}_i)} [f_k(\mathbf{x}_i, \mathbf{y}')] - f_k(\mathbf{x}_i, \mathbf{y}_i)}_{\text{---}}$$



↑
Hmm... there are exponentially many possible \mathbf{y}' here,
we shall find a way to calculate this term efficiently!

Features for CRF

$$\mathbf{E}_{p(\mathbf{y}'|\mathbf{x}_i)} [f_k(\mathbf{x}_i, \mathbf{y}')] \quad \text{💬}$$

Let's consider a
very simple type
of feature

$$f_{918}(\mathbf{x}_i, \mathbf{y}') = \sum_j \llbracket j\text{-th tag is "N" and } (j+1)\text{-th tag is "V"} \rrbracket$$

$$= 1$$

\mathbf{y}'	N	A	N	V	V
\mathbf{x}_i	<i>Fruit</i>	<i>flies</i>	<i>like</i>	<i>a</i>	<i>banana</i>

Features for CRF

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

Let's consider a
very simple type
of feature

$$f_{918}(\mathbf{x}_i, \mathbf{y}') = \sum_j \llbracket j\text{-th tag is "N" and } (j+1)\text{-th tag is "V"} \rrbracket$$

$$= 2$$



Features for CRF

$$\mathbf{E}_{p(\mathbf{y}'|x_i)}[f_k(\mathbf{x}_i, \mathbf{y}')]$$

						$p(\mathbf{y}' x_i)$	$f_{918}(\mathbf{x}_i, \mathbf{y}')$
\mathbf{y}'	A	A	A	A	A	0.001	0
	...						
	A	N	V	N	V	0.007	2
	A	N	V	D	N	0.002	1
\mathbf{x}_i	Fruit	flies	like	a	banana		



Features for CRF

$$\mathbf{E}_{p(\mathbf{y}'|x_i)} [f_k(\mathbf{x}_i, \mathbf{y}')]]$$

						$p(\mathbf{y}' x_i)$	$f_{918}(\mathbf{x}_i, \mathbf{y}')$
\mathbf{y}'	A	A	A	A	A	0.001	0
	...						
	A	N	V	N	V	0.007	2
	A	N	V	D	N	0.002	1
\mathbf{x}_i	Fruit	flies	like	a	banana		

However, it is not feasible to enumerate all the \mathbf{y}' sequences. Why?

Features for CRF

$$\mathbf{E}_{p(\mathbf{y}'|x_i)} [f_k(\mathbf{x}_i, \mathbf{y}')]]$$

						$p(\mathbf{y}' x_i)$	$f_{918}(\mathbf{x}_i, \mathbf{y}')$
\mathbf{y}'	A	A	A	A	A	0.001	0
	...						
	A	N	V	N	V	0.007	2
	A	N	V	D	N	0.002	1
\mathbf{x}_i	Fruit	flies	like	a	banana		



Hold on a second... isn't this the same as what we did in HMM for calculating the expected counts for transitions?

Features for CRF

$$\mathbf{E}_{p(\mathbf{y}'|x_i)} [f_k(\mathbf{x}_i, \mathbf{y}')]]$$

						$p(\mathbf{y}' x_i)$	$f_{918}(\mathbf{x}_i, \mathbf{y}')$
\mathbf{y}'	A	A	A	A	A	0.001	0
	...						
	A	N	V	N	V	0.007	2
	A	N	V	D	N	0.002	1
\mathbf{x}_i	Fruit	flies	like	a	banana		



Hold on a second... isn't this the same as what we did with the
HMM for calculating the expected co-occurrences?

Forward-backward
algorithm!

Features for CRF

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

Another type of
simple features

$$f_{102}(\mathbf{x}_i, \mathbf{y}') = \sum_j \llbracket j\text{-th word is "a" and } j\text{-th tag is "D"} \rrbracket$$

$$= 1$$



Features for CRF

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

Another type of
simple features

$$f_{102}(\mathbf{x}_i, \mathbf{y}') = \sum_j \llbracket j\text{-th word is "a" and } j\text{-th tag is "D"} \rrbracket$$

$$= 0$$

\mathbf{y}'	A	A	A	A	A
\mathbf{x}_i	<i>Fruit</i>	<i>flies</i>	<i>like</i>	<i>a</i>	<i>banana</i>

Features for CRF

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

Yet another
feature



$$f_{1782}(\mathbf{x}_i, \mathbf{y}') = \sum_j \llbracket j\text{-th word is "flies" and } j\text{-th tag is "N" and } (j + 1)\text{-th tag is "V"} \rrbracket$$

$$= 1$$

\mathbf{y}' **A**
 \mathbf{x}_i *Fruit*

N **V**
flies *like*

D
a

N
banana

Features for CRF

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

Yet another
feature

$$f_{1782}(\mathbf{x}_i, \mathbf{y}') = \sum_j \llbracket j\text{-th word is "flies" and } j\text{-th tag is "N" and } (j+1)\text{-th tag is "V"} \rrbracket$$

$$= 0$$

\mathbf{y}'	N	A	N	V	V
\mathbf{x}_i	<i>Fruit</i>	<i>flies</i>	<i>like</i>	<i>a</i>	<i>banana</i>

Features for CRF

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

$$f_{102}(\mathbf{x}_i, \mathbf{y}') =$$

$$\sum_j \llbracket j\text{-th word is "a" and } j\text{-th tag is "D"} \rrbracket$$

$$f_{918}(\mathbf{x}_i, \mathbf{y}') =$$

$$\sum_j \llbracket j\text{-th tag is "N" and } (j+1)\text{-th tag is "V"} \rrbracket$$

$$f_{1782}(\mathbf{x}_i, \mathbf{y}') =$$

$$\sum_j \llbracket j\text{-th word is "flies" and } j\text{-th tag is "N" and } (j+1)\text{-th tag is "V"} \rrbracket$$

\mathbf{y}'	N	A	N	V	V
\mathbf{x}_i	<i>Fruit</i>	<i>flies</i>	<i>like</i>	<i>a</i>	<i>banana</i>



We shall define *local* features so as to apply efficient dynamic programming algorithms for calculating expected feature counts!

Learning in CRF

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp(\mathbf{f}(\mathbf{x}, \mathbf{y}) \cdot \boldsymbol{\theta})}{\sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}, \mathbf{y}') \cdot \boldsymbol{\theta})}$$

$$\min_{\boldsymbol{\theta}} - \sum_i \log p(\mathbf{y}_i \mid \mathbf{x}_i)$$

$$\frac{\partial -\log p(\mathbf{y}_i \mid \mathbf{x}_i)}{\partial \theta_k} = \underbrace{\mathbf{E}_{p(\mathbf{y}' \mid \mathbf{x}_i)} [f_k(\mathbf{x}_i, \mathbf{y}')] - f_k(\mathbf{x}_i, \mathbf{y}_i)}_{\text{Forward-backward algorithm}}$$

Forward-backward algorithm

Decoding in CRF

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{\exp(\mathbf{f}(\mathbf{x}, \mathbf{y}) \cdot \boldsymbol{\theta})}{\sum_{\mathbf{y}'} \exp(\mathbf{f}(\mathbf{x}, \mathbf{y}') \cdot \boldsymbol{\theta})}$$

$$\begin{aligned} \mathbf{y}^* &= \arg \max_{\mathbf{y}} \log p(\mathbf{y} \mid \mathbf{x}^{(new)}) \\ &= \arg \max_{\mathbf{y}} \mathbf{f}(\mathbf{x}^{(new)}, \mathbf{y}) \cdot \boldsymbol{\theta} \end{aligned}$$

A new input
sentence

Learned model
parameters



Viterbi algorithm!

CRF

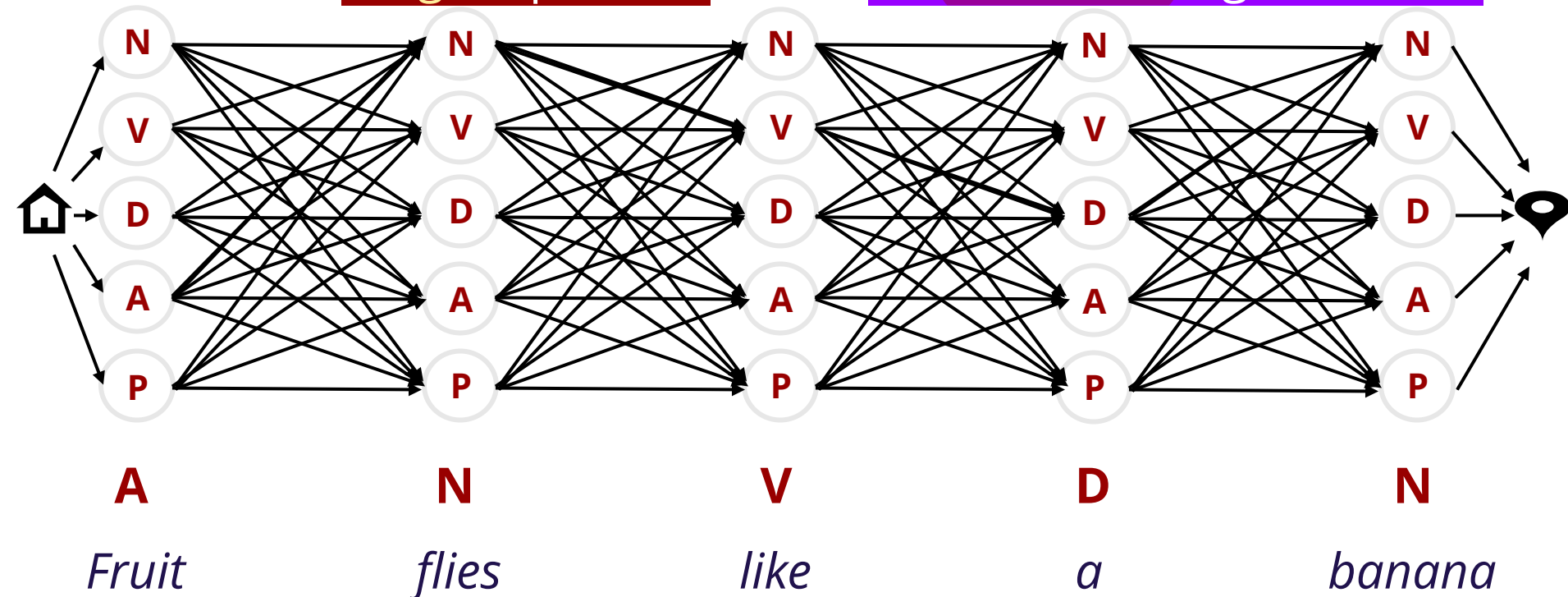
Objective Function



$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \boldsymbol{\theta} + \log \sum_y \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \boldsymbol{\theta}) \right)$$

Score of the
gold path

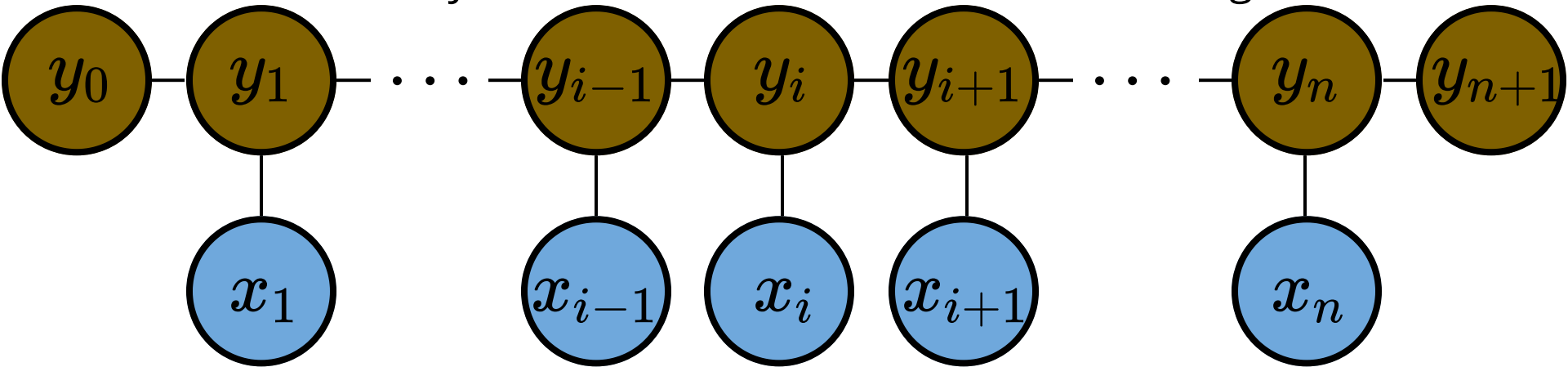
Calculated with forward-
backward algorithm



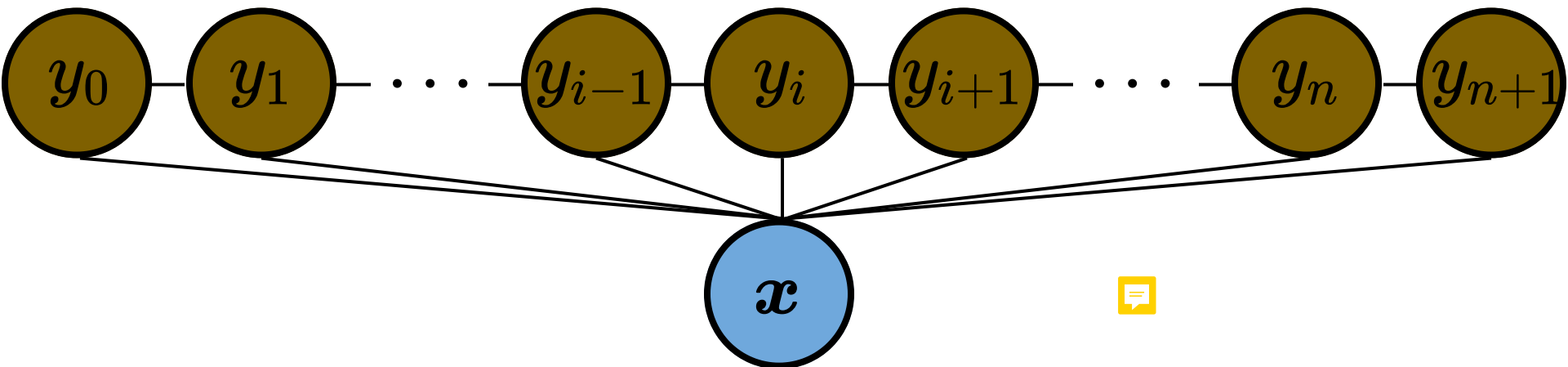
Optional

CRF

It is OK for you to think of the model as something like:



Although it actually looks something like this:



Discriminative Models

Probabilistic vs Non-Probabilistic




Probabilistic	Non-Probabilistic
Logistic regression	Perceptron
Softmax regression	Support vector machines



Is it possible to train a discriminative model for sequence labeling with a non-probabilistic model such as Perceptron?

Perceptron

$$\hat{y} \in \{-1, +1\}$$


$$\hat{y} \leftarrow \text{sign}(\mathbf{x}_i \cdot \boldsymbol{\theta} + \theta_0)$$

if $\hat{y} \neq y_i$ then

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y_i \mathbf{x}_i$$

$$\theta_0 \leftarrow \theta_0 + y_i$$

Invented more than 50 years ago.
Still popular today due to its simplicity!

Perceptron

Let's understand it a little better...

$$\hat{y} \leftarrow \arg \max_y y(\mathbf{x}_i \cdot \boldsymbol{\theta} + \theta_0)$$



$$\hat{y} \leftarrow \text{sign}(\mathbf{x}_i \cdot \boldsymbol{\theta} + \theta_0)$$

if $\hat{y} \neq y_i$ then

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y_i \mathbf{x}_i$$

$$\theta_0 \leftarrow \theta_0 + y_i$$

Perceptron

$$\hat{y} \leftarrow \arg \max_y y(\mathbf{x}_i \cdot \boldsymbol{\theta} + \theta_0)$$

if $\hat{y} \neq y_i$ then

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + y_i \mathbf{x}_i$$

$$\theta_0 \leftarrow \theta_0 + y_i$$

Perceptron

$$\hat{y} \leftarrow \arg \max_y \begin{pmatrix} yx_i/2 \\ y/2 \end{pmatrix} \cdot \begin{pmatrix} \theta \\ \theta_0 \end{pmatrix}$$



$$\hat{y} \leftarrow \arg \max_y y(x_i \cdot \theta + \theta_0)$$



if $\hat{y} \neq y_i$ then

$$\theta \leftarrow \theta + y_i x_i$$

$$\theta_0 \leftarrow \theta_0 + y_i$$



$$\begin{pmatrix} \theta \\ \theta_0 \end{pmatrix} \leftarrow \begin{pmatrix} \theta \\ \theta_0 \end{pmatrix} + \begin{pmatrix} y_i x_i \\ y_i \end{pmatrix}$$

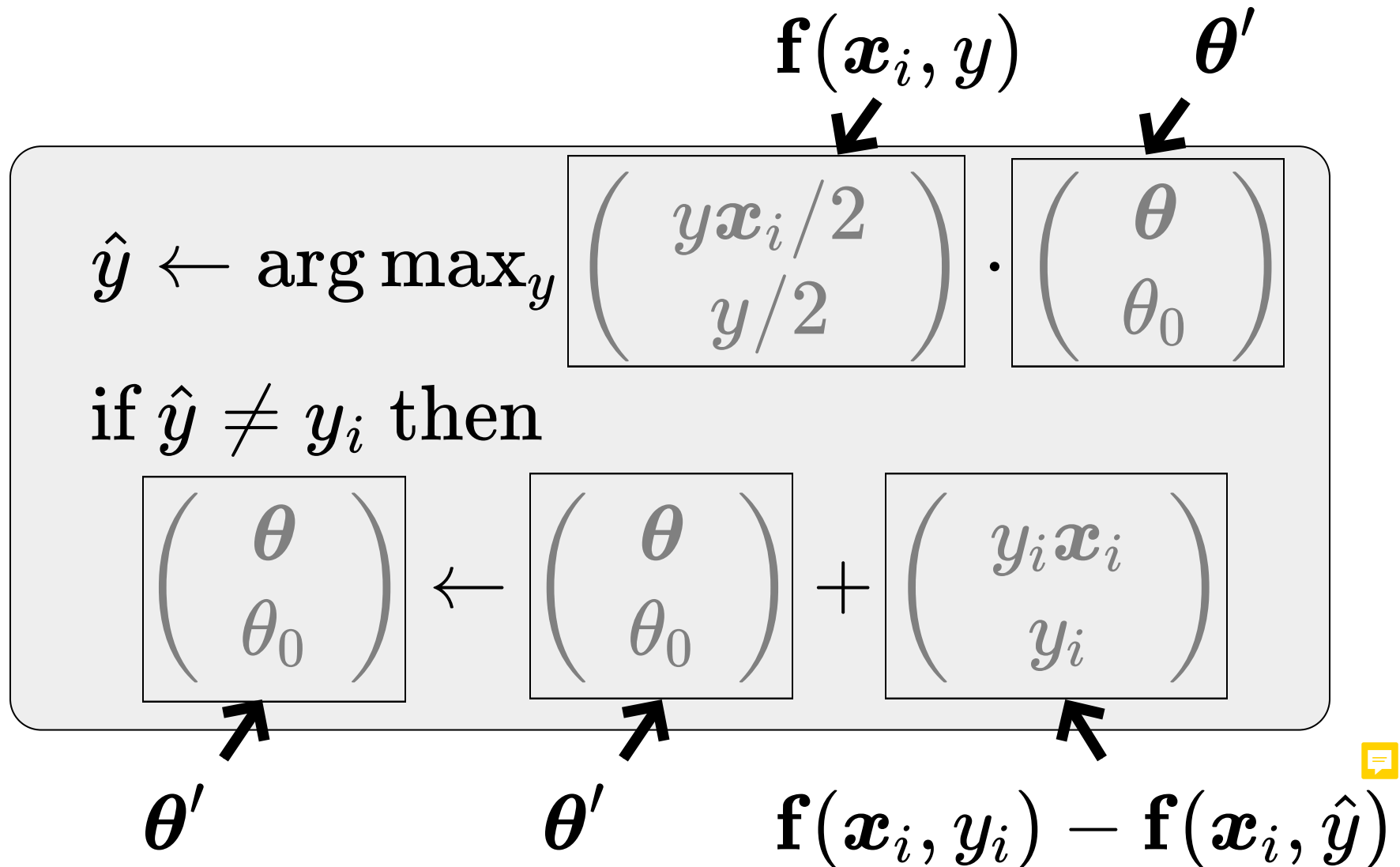
Perceptron

$$\hat{y} \leftarrow \arg \max_y \begin{pmatrix} y\mathbf{x}_i/2 \\ y/2 \end{pmatrix} \cdot \begin{pmatrix} \boldsymbol{\theta} \\ \theta_0 \end{pmatrix}$$

if $\hat{y} \neq y_i$ then

$$\begin{pmatrix} \boldsymbol{\theta} \\ \theta_0 \end{pmatrix} \leftarrow \begin{pmatrix} \boldsymbol{\theta} \\ \theta_0 \end{pmatrix} + \begin{pmatrix} y_i\mathbf{x}_i \\ y_i \end{pmatrix}$$

Perceptron



Perceptron

$$\hat{y} \leftarrow \arg \max_y \mathbf{f}(\mathbf{x}_i, y) \cdot \boldsymbol{\theta}'$$

if $\hat{y} \neq y_i$ then

$$\boldsymbol{\theta}' \leftarrow \boldsymbol{\theta}' + \mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})$$

This allows us to make a generalization for the variable y : from a binary output to a structured output

Structured Perceptron

(Collins, 1999)

This is a structure now!

The space of structures given x_i

$$\hat{y}_i \leftarrow \arg \max_{y \in \text{GEN}(x_i)} \mathbf{f}(x_i, y) \cdot \theta$$

if $\hat{y}_i \neq y_i$ then

$$\theta \leftarrow \theta + \mathbf{f}(x_i, y_i) - \mathbf{f}(x_i, \hat{y}_i)$$

Structured Perceptron

$$\hat{\mathbf{y}}_i \leftarrow \arg \max_{\mathbf{y} \in \mathbf{GEN}(\mathbf{x}_i)} \mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \boldsymbol{\theta}$$

if $\hat{\mathbf{y}}_i \neq \mathbf{y}_i$ then

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{f}(\mathbf{x}_i, \hat{\mathbf{y}}_i)$$

We shall make sure searching in the space defined by the $\mathbf{GEN}(\mathbf{x}_i)$ function is efficient.



Viterbi algorithm!
Define similar **local features**
as what we did in CRF!

Structured Perceptron

$$\hat{\mathbf{y}}_i \leftarrow \arg \max_{\mathbf{y} \in \text{GEN}(\mathbf{x}_i)} \mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \boldsymbol{\theta}$$

if $\hat{\mathbf{y}}_i \neq \mathbf{y}_i$ then

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) - \mathbf{f}(\mathbf{x}_i, \hat{\mathbf{y}}_i)$$

This update rule is essentially stochastic gradient descent
with learning rate 1!

What is the underlying objective function?

$$- \mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \boldsymbol{\theta} + \max_{\mathbf{y}} (\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \boldsymbol{\theta})$$

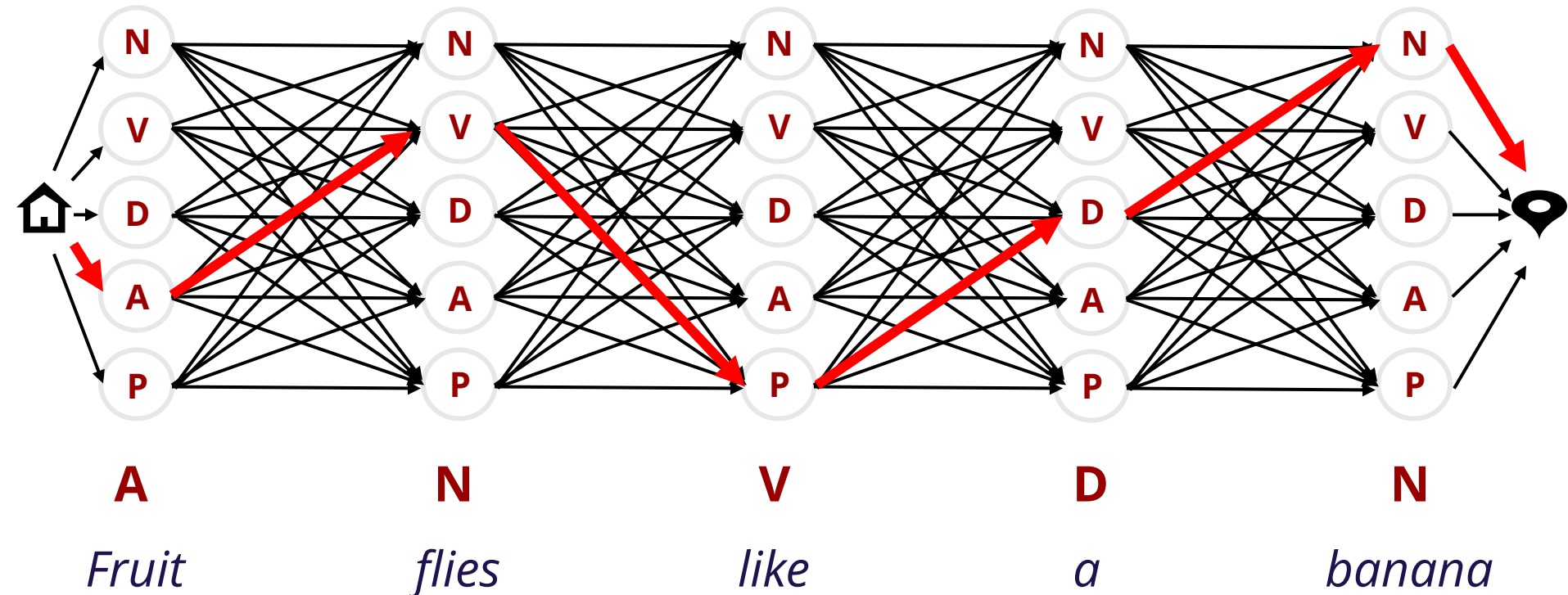
Structured Perceptron

Objective Function

$$\min_{\theta} \sum_i \left(- \mathbf{f}(x_i, y_i) \cdot \theta + \max_y \left(\mathbf{f}(x_i, y) \cdot \theta \right) \right)$$

Score of the
gold path

Score of the
best path

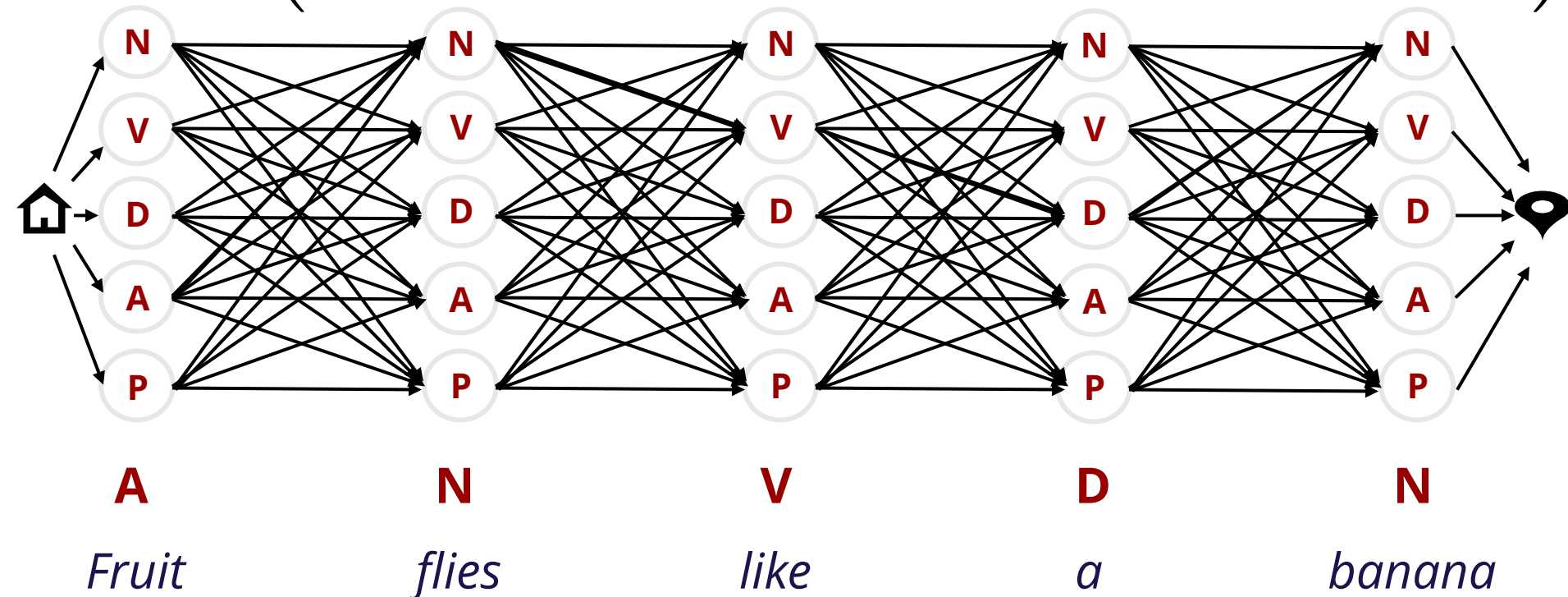


CRF vs SP

Comparison

$$\min_{\theta} \sum_i \left(-\mathbf{f}(x_i, y_i) \cdot \theta + \log \sum_y \exp(\mathbf{f}(x_i, y) \cdot \theta) \right)$$

$$\min_{\theta} \sum_i \left(-\mathbf{f}(x_i, y_i) \cdot \theta + \max_y (\mathbf{f}(x_i, y) \cdot \theta) \right)$$



CRF vs SP

Comparison

$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \log \sum_y \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta) \right)$$
$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \max_y (\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta) \right)$$

The difference!

CRF vs SP

Comparison

$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \log \sum_y \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta) \right)$$
$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \max_y (\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta) \right)$$

$$\mathcal{Y} = \{1, 3, 4, 8, 18\}$$

$$\log \sum_{y \in \mathcal{Y}} \exp(y) = ?$$

$$\max_{y \in \mathcal{Y}} (y) = ?$$



CRF vs SP

Comparison

$$\begin{aligned} \min_{\theta} \sum_i & \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \log \sum_y \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta) \right) \\ \min_{\theta} \sum_i & \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \max_y (\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta) \right) \end{aligned}$$

$$\mathcal{Y} = \{1, 3, 4, 8, 18\}$$

$$\log \sum_{y \in \mathcal{Y}} \exp(y) = 18.0000465777..$$

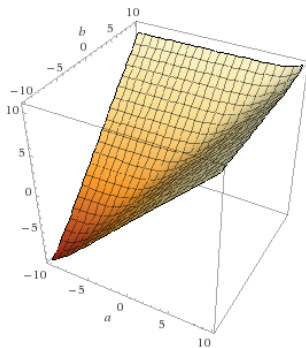
$$\max_{y \in \mathcal{Y}} (y) = 18$$

CRF vs SP

Comparison

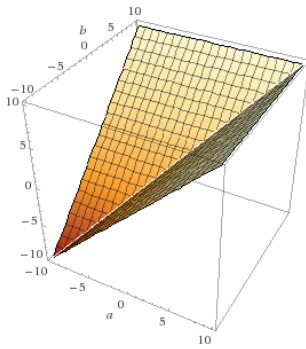
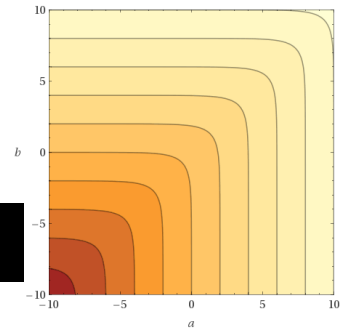
$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \log \sum_y \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta) \right)$$

$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \max_y (\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta) \right)$$



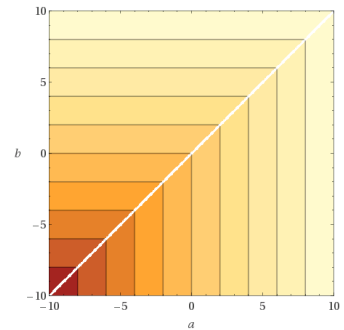
Computed by Wolfram|Alpha

$\log(\exp(a) + \exp(b))$
A "soft"/"smooth" version of max!



Computed by Wolfram|Alpha

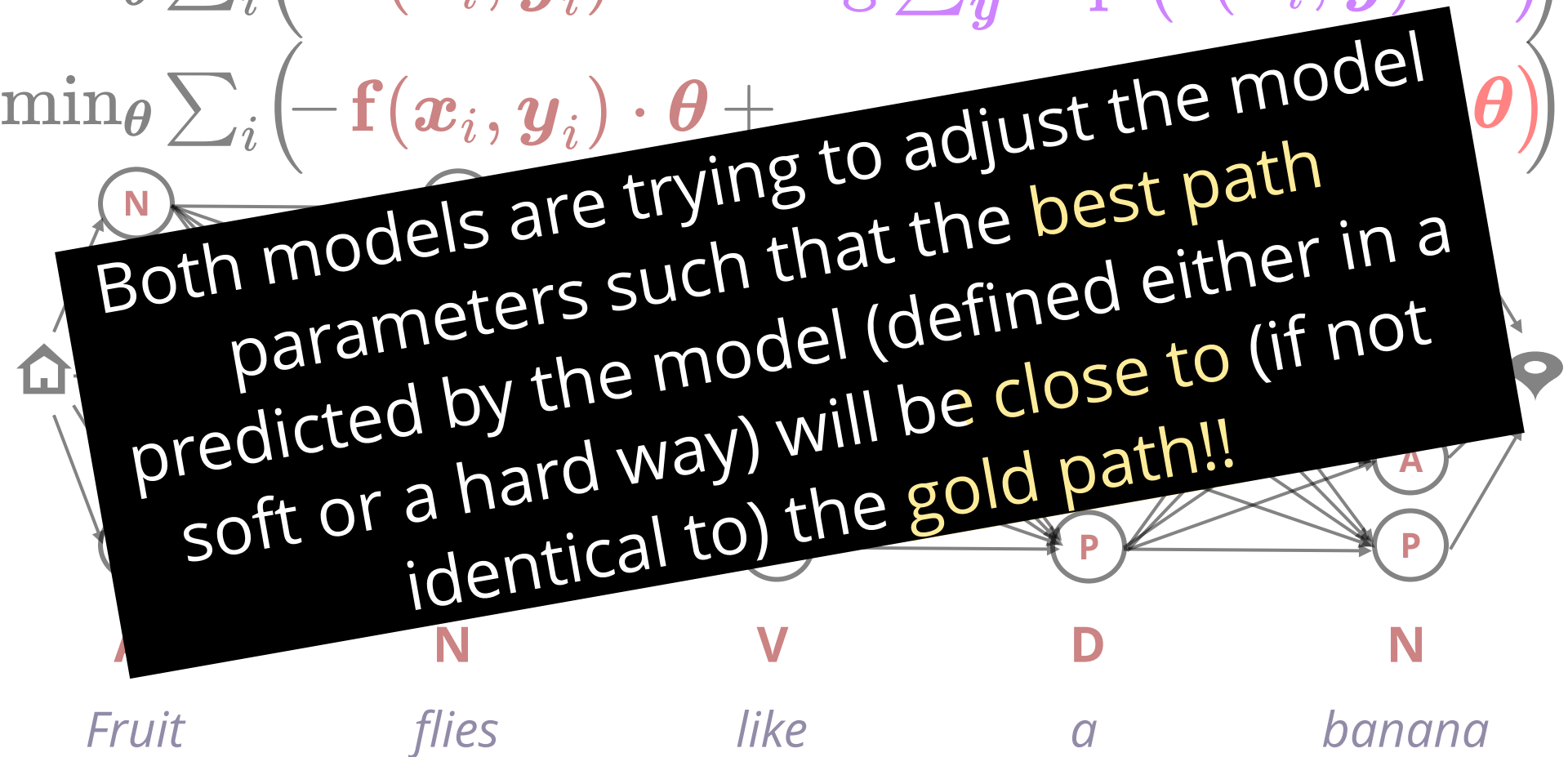
$\max(a, b)$



CRF vs SP

Comparison

$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \boldsymbol{\theta} + \log \sum_y \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \boldsymbol{\theta}) \right)$$
$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \boldsymbol{\theta} + \log \sum_y \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \boldsymbol{\theta}) \right)$$



Sequence Labeling

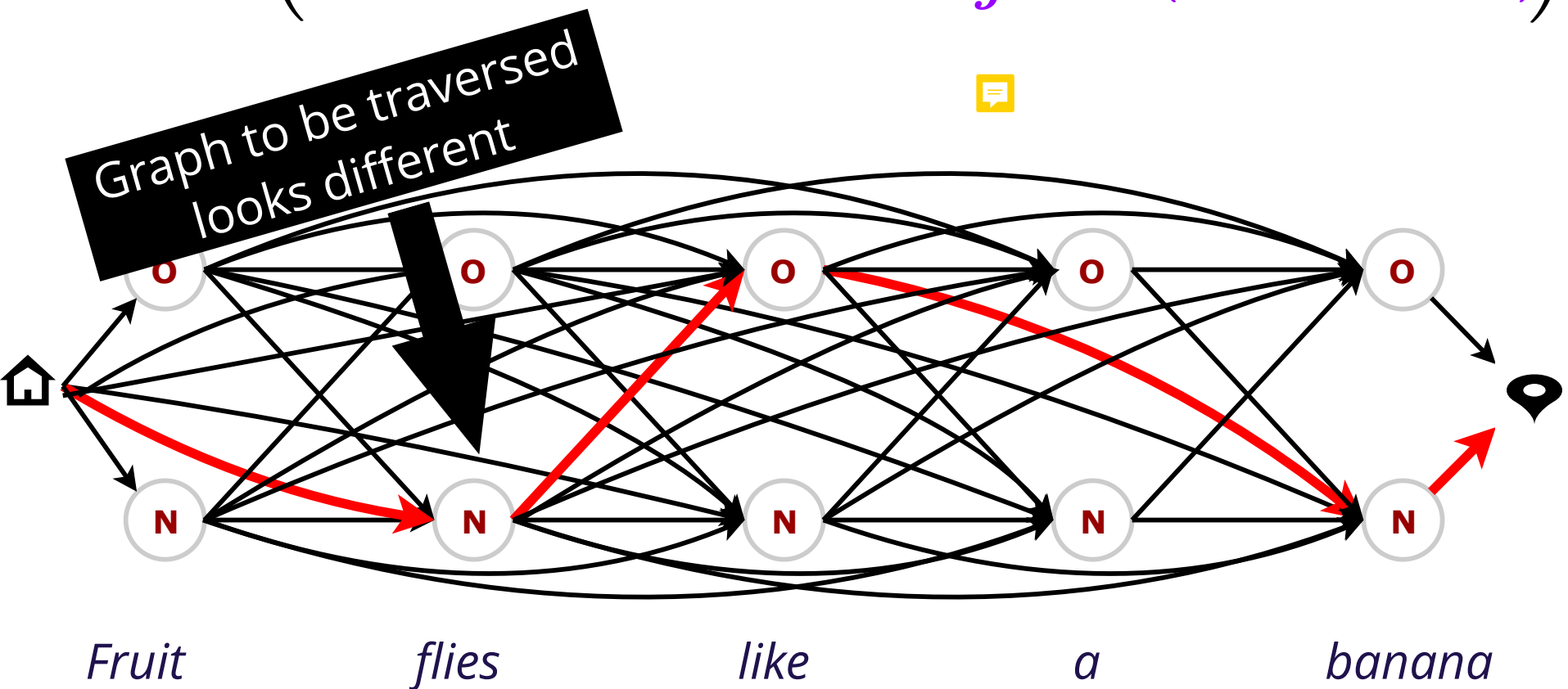


Model	Pros	Cons
Hidden Markov Model (HMM)	Probabilistic, efficient training	Unable to exploit linguistic features
Maximum Entropy Markov Model (MEMM)	Probabilistic, able to exploit features	Training is slower than HMM, local normalization
Conditional Random Fields (CRF)	Probabilistic, able to exploit features, global normalization	Training may be slower than MEMM
Structured Perceptron (SP)	Able to exploit features, performance comparable to CRF	Non-probabilistic

Optional

Semi-Markov CRF (Sarawagi & Cohen, 2005)

$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \boldsymbol{\theta} + \log \sum_y \exp(\mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \boldsymbol{\theta}) \right)$$

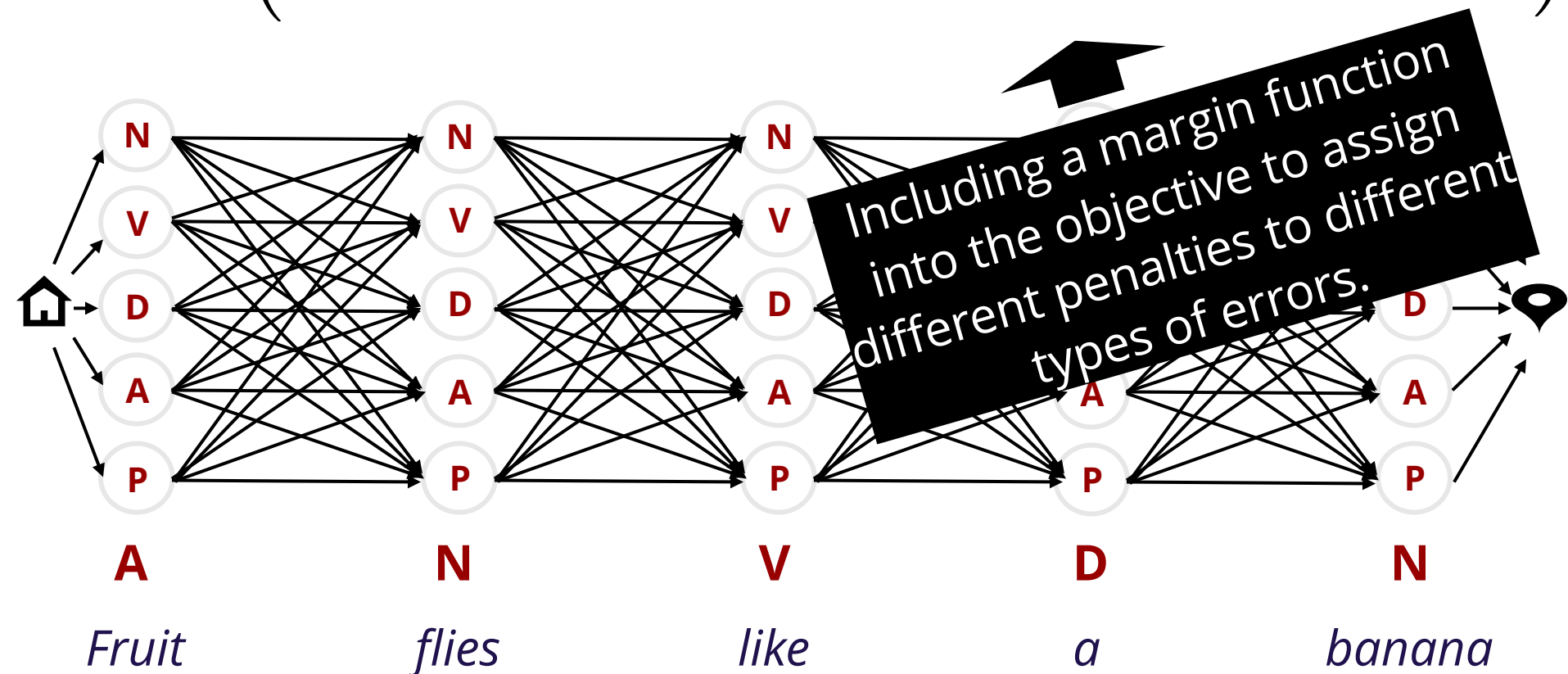


Optional

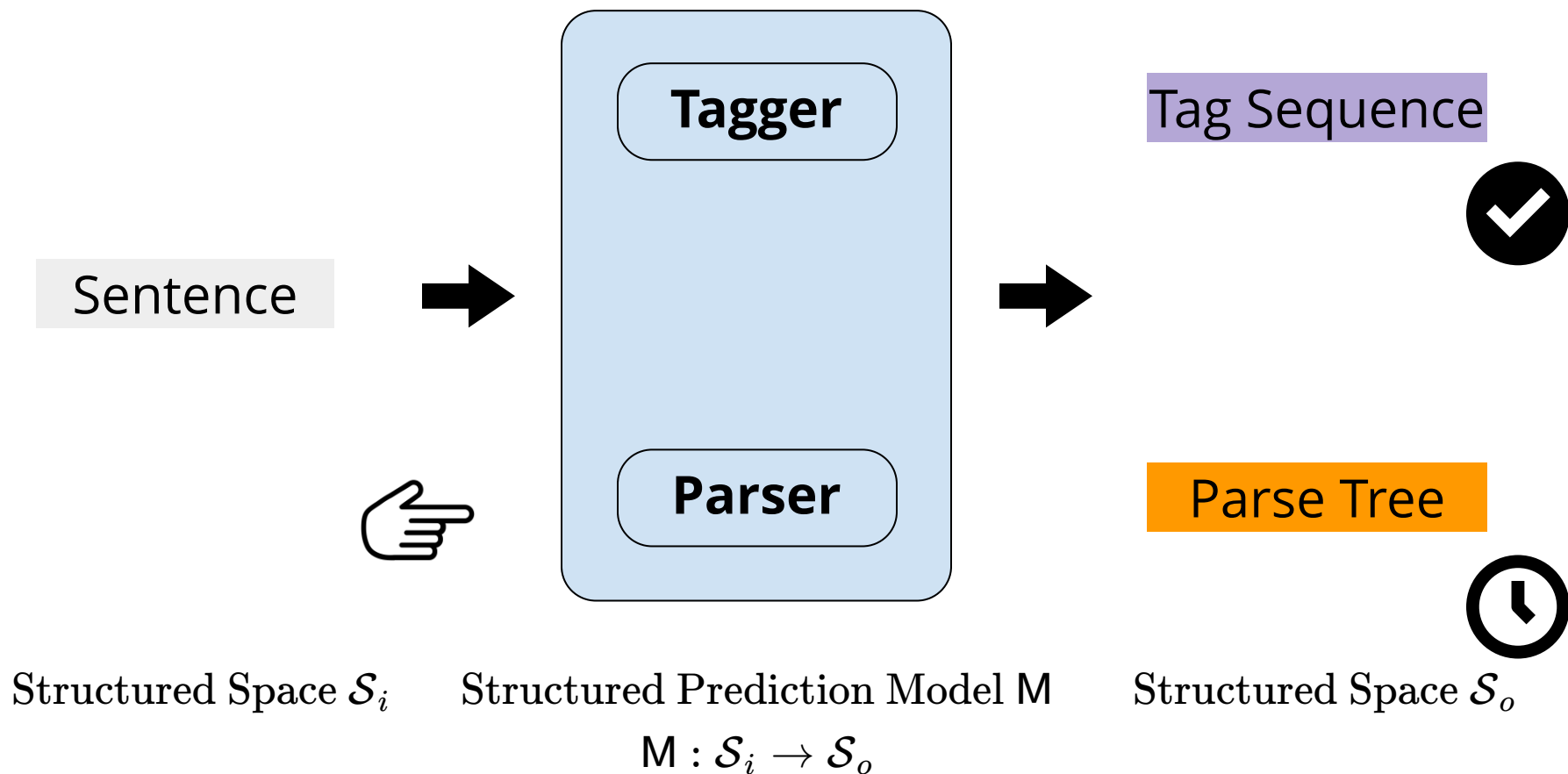
Structural SVM

(Tsochantaridis et al., 2005)

$$\min_{\theta} \sum_i \left(-\mathbf{f}(\mathbf{x}_i, \mathbf{y}_i) \cdot \theta + \max_y \left(\Delta(\mathbf{y}_i, \mathbf{y}) + \mathbf{f}(\mathbf{x}_i, \mathbf{y}) \cdot \theta \right) \right)$$



Structured Prediction



Tasks in NLP

