

Sequential Data Modeling - The Structured Perceptron

Graham Neubig
Nara Institute of Science and Technology (NAIST)

Prediction Problems

Given x , predict y

Prediction Problems

Given x , predict y

A book review

Oh, man I love this book!
This book is so boring...

Is it positive?

yes
no

Binary
Prediction
(2 choices)

A tweet

On the way to the park!
公園に行くなう！

Its language

English
Japanese

Multi-class
Prediction
(several choices)

A sentence

I read a book

Its parts-of-speech

N	VBD	DET	NN
I	read	a	book

Structured
Prediction
(millions of choices)

Sequential prediction is a subset

Simple Prediction: The Perceptron Model

Example we will use:

- Given an introductory sentence from Wikipedia
- Predict **whether the article is about a person**

<u>Given</u>		<u>Predict</u>
Gonso was a Sanron sect priest (754-827) in the late Nara and early Heian periods.	→	Yes!
Shichikuzan Chigogataki Fudomyoo is a historical site located at Magura, Maizuru City, Kyoto Prefecture.	→	No!

- This is **binary classification** (of course!)

How do We Predict?

Gonso was a Sanron sect priest (754 – 827)
in the late Nara and early Heian periods .

Shichikuzan Chigogataki Fudomyoo is
a historical site located at Magura , Maizuru
City , Kyoto Prefecture .

How do We Predict?

Contains “priest” →
probably person!

Contains “(<#>-<#>)” →
probably person!

Gonso was a Sanron sect **priest** (**754 – 827**)
in the late Nara and early Heian periods .

Contains
“site” →
probably not person!

Shichikuzan Chigogataki Fudomyoo is
a historical **site** located at Magura , Maizuru
City , **Kyoto Prefecture** .

Contains
“Kyoto Prefecture” →
probably not person!

Combining Pieces of Information

- Each element that helps us predict is a *feature*

contains “priest”	contains “(<#>-<#>)”
contains “site”	contains “Kyoto Prefecture”

- Each feature has a weight, *positive* if it indicates “yes”, and *negative* if it indicates “no”

$W_{\text{contains “priest”}} = 2$	$W_{\text{contains “(<#>-<#>)”}} = 1$
$W_{\text{contains “site”}} = -3$	$W_{\text{contains “Kyoto Prefecture”}} = -1$

- For a new example, sum the weights

Kuya (903-972) was a priest born in Kyoto Prefecture.	$2 + -1 + 1 = 2$
--	------------------

- If the sum is at least 0: “yes”, otherwise: “no”

Let me Say that in Math!

$$\begin{aligned} y &= \text{sign}(\mathbf{w} \cdot \boldsymbol{\varphi}(\mathbf{x})) \\ &= \text{sign}\left(\sum_{i=1}^I w_i \cdot \varphi_i(\mathbf{x})\right) \end{aligned}$$

- \mathbf{x} : the input
- $\boldsymbol{\varphi}(\mathbf{x})$: vector of feature functions $\{\varphi_1(\mathbf{x}), \varphi_2(\mathbf{x}), \dots, \varphi_I(\mathbf{x})\}$
- \mathbf{w} : the weight vector $\{w_1, w_2, \dots, w_I\}$
- y : the prediction, +1 if “yes”, -1 if “no”
 - ($\text{sign}(v)$ is +1 if $v \geq 0$, -1 otherwise)

Example Feature Functions: Unigram Features

- Equal to “number of times a particular word appears”

$x =$ A site , located in Maizuru , Kyoto

$$\varphi_{\text{unigram "A"}}(x) = 1 \quad \varphi_{\text{unigram "site"}}(x) = 1 \quad \varphi_{\text{unigram ","}}(x) = 2$$

$$\varphi_{\text{unigram "located"}}(x) = 1 \quad \varphi_{\text{unigram "in"}}(x) = 1$$

$$\varphi_{\text{unigram "Maizuru"}}(x) = 1 \quad \varphi_{\text{unigram "Kyoto"}}(x) = 1$$

$$\left. \begin{array}{l} \varphi_{\text{unigram "the"}}(x) = 0 \quad \varphi_{\text{unigram "temple"}}(x) = 0 \\ \dots \end{array} \right\} \text{The rest are all 0}$$

- For convenience, we use feature names ($\varphi_{\text{unigram "A"}}$) instead of feature indexes (φ_1)

Calculating the Weighted Sum

$x =$ A site , located in Maizuru , Kyoto

$\varphi_{\text{unigram "A"}}(x) = 1$	$W_{\text{unigram "a"}} = 0$	0	+
$\varphi_{\text{unigram "site"}}(x) = 1$	$W_{\text{unigram "site"}} = -3$	-3	+
$\varphi_{\text{unigram "located"}}(x) = 1$	$W_{\text{unigram "located"}} = 0$	0	+
$\varphi_{\text{unigram "Maizuru"}}(x) = 1$	$W_{\text{unigram "Maizuru"}} = 0$	0	+
$\varphi_{\text{unigram " , "}}(x) = 2$	$W_{\text{unigram " , "}} = 0$	0	+
$\varphi_{\text{unigram "in"}}(x) = 1$	$W_{\text{unigram "in"}} = 0$	0	+
$\varphi_{\text{unigram "Kyoto"}}(x) = 1$	$W_{\text{unigram "Kyoto"}} = 0$	0	+
$\varphi_{\text{unigram "priest"}}(x) = 0$	$W_{\text{unigram "priest"}} = 2$	0	+
$\varphi_{\text{unigram "black"}}(x) = 0$	$W_{\text{unigram "black"}} = 0$	0	+
\dots	\dots	$=$	
		-3	\rightarrow No!

Learning Weights Using the Perceptron Algorithm

Learning Weights

- **Manually creating weights is hard**
 - Many many possible useful features
 - Changing weights changes results in unexpected ways
- **Instead, we can learn from labeled data**

y	x
1	FUJIWARA no Chikamori (year of birth and death unknown) was a samurai and poet who lived at the end of the Heian period .
1	Ryonen (1646 - October 29 , 1711) was a Buddhist nun of the Obaku Sect who lived from the early Edo period to the mid-Edo period .
-1	A moat settlement is a village surrounded by a moat .
-1	Fushimi Momoyama Athletic Park is located in Momoyama-cho , Kyoto City , Kyoto Prefecture .

Online Learning

```
create map  $w$ 
for / iterations
  for each labeled pair  $x, y$  in the data
     $\phi = \text{CREATE\_FEATURES}(x)$ 
     $y' = \text{PREDICT\_ONE}(w, \phi)$ 
    if  $y' \neq y$ 
       $\text{UPDATE\_WEIGHTS}(w, \phi, y)$ 
```

- In other words
 - Try to classify each training example
 - Every time we make a mistake, update the weights
- Many different online learning algorithms
 - The most simple is the **perceptron**

Perceptron Weight Update

$$\mathbf{w} \leftarrow \mathbf{w} + y \boldsymbol{\varphi}(\mathbf{x})$$

- In other words:
 - If $y=1$, increase the weights for features in $\boldsymbol{\varphi}(\mathbf{x})$
 - Features for positive examples get a higher weight
 - If $y=-1$, decrease the weights for features in $\boldsymbol{\varphi}(\mathbf{x})$
 - Features for negative examples get a lower weight
- Every time we update, our predictions get better!

Example: Initial Update

- Initialize $\mathbf{w}=0$

\mathbf{x} = A site , located in Maizuru , Kyoto $y = -1$

$$\mathbf{w} \cdot \boldsymbol{\varphi}(\mathbf{x}) = 0 \qquad y' = \text{sign}(\mathbf{w} \cdot \boldsymbol{\varphi}(\mathbf{x})) = 1$$

$$y' \neq y$$

$$\mathbf{w} \leftarrow \mathbf{w} + y \boldsymbol{\varphi}(\mathbf{x})$$

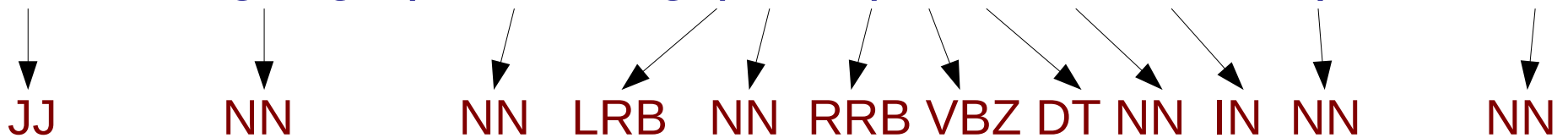
$W_{\text{unigram "Maizuru"}} = -1$	$W_{\text{unigram "A"}} = -1$
$W_{\text{unigram ","}} = -2$	$W_{\text{unigram "site"}} = -1$
$W_{\text{unigram "in"}} = -1$	$W_{\text{unigram "located"}} = -1$
$W_{\text{unigram "Kyoto"}} = -1$	

Review: The HMM Model

Probabilistic Model for Tagging

- “Find the most probable **tag sequence**, given **the sentence**”

Natural language processing (NLP) is a field of computer science



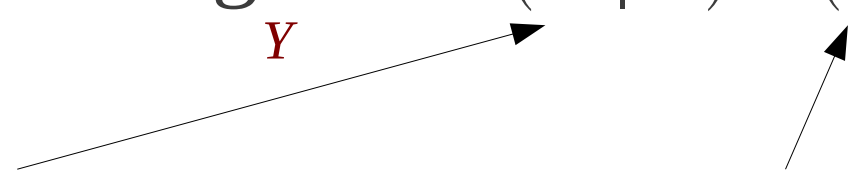
$$\operatorname{argmax}_Y P(Y|X)$$

- Any ideas?

Generative Sequence Model

- First decompose probability using **Bayes' law**

$$\operatorname{argmax}_Y P(Y|X) = \operatorname{argmax}_Y \frac{P(X|Y)P(Y)}{P(X)}$$

$$= \operatorname{argmax}_Y P(X|Y)P(Y)$$


Model of word/POS interactions
“natural” is probably a JJ

Model of POS/POS interactions
NN comes after DET

- Also sometimes called the **“noisy-channel model”**

Hidden Markov Models (HMMs) for POS Tagging

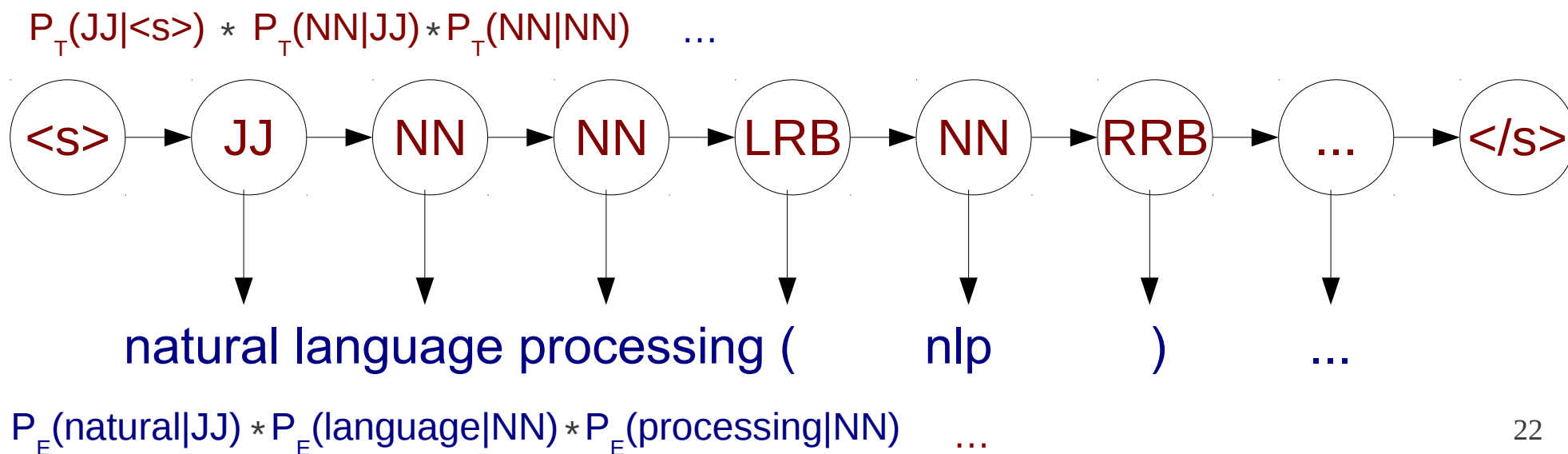
- POS → POS **transition** probabilities

- Like a bigram model!

$$P(Y) \approx \prod_{i=1}^{l+1} P_T(y_i | y_{i-1})$$

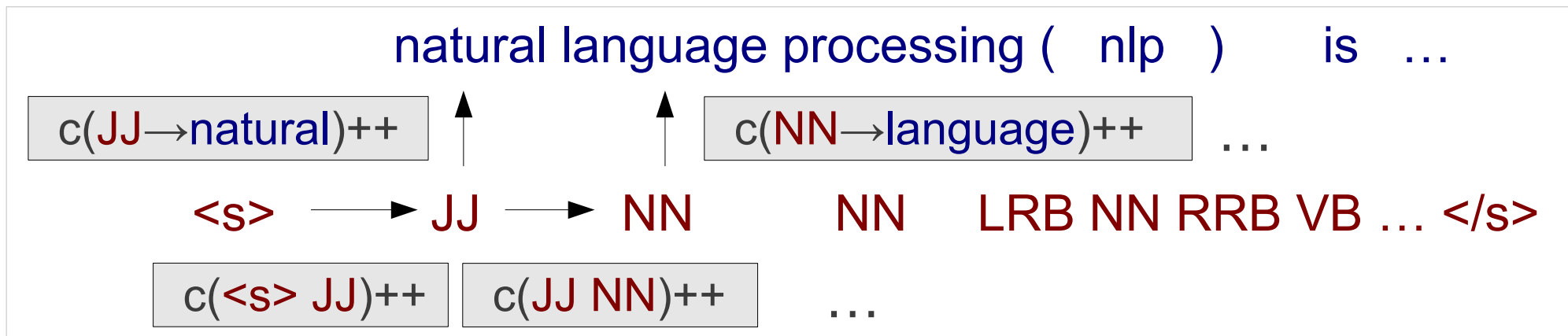
- POS → Word **emission** probabilities

$$P(X|Y) \approx \prod_{i=1}^l P_E(x_i | y_i)$$



Learning Markov Models (with tags)

- Count the number of occurrences in the corpus and



- Divide by context to get probability

$$P_T(\text{LRB}|\text{NN}) = c(\text{NN LRB})/c(\text{NN}) = 1/3$$

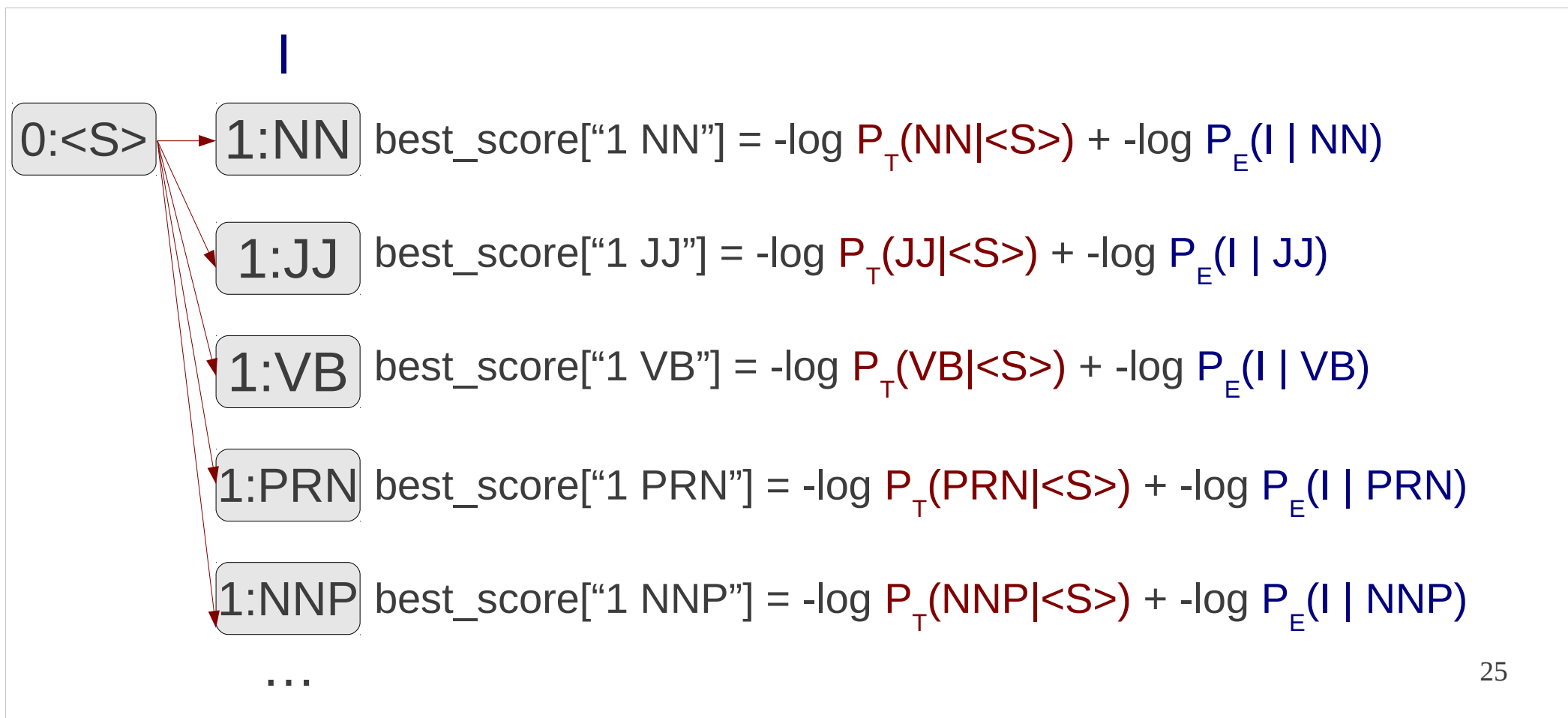
$$P_E(\text{language}|\text{NN}) = c(\text{NN} \rightarrow \text{language})/c(\text{NN}) = 1/3$$

Remember: HMM Viterbi Algorithm

- **Forward step**, calculate the best path to a node
 - Find the path to each node with the **lowest negative log probability**
- **Backward step**, reproduce the path
 - This is easy, almost the same as word segmentation

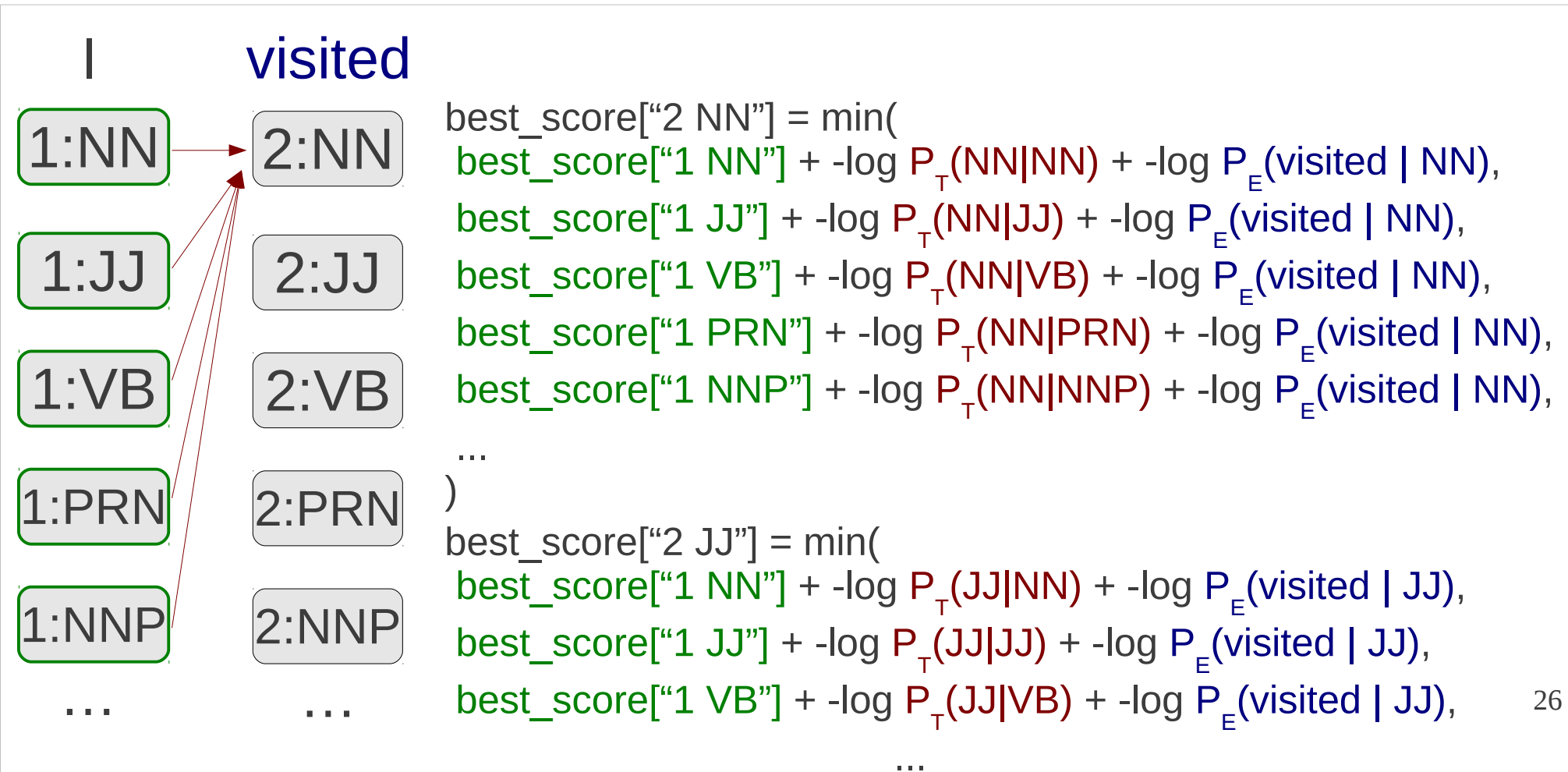
Forward Step: Part 1

- First, calculate transition from $\langle S \rangle$ and emission of the first word for every POS



Forward Step: Middle Parts

- For middle words, calculate the minimum score for all possible previous POS tags



The Structured Perceptron

So Far, We Have Learned

Classifiers

Perceptron

Lots of features

Binary prediction

Generative Models

HMM

Conditional probabilities

Structured prediction

Structured Perceptron

Classifiers

Perceptron

Lots of features

Binary prediction

Generative Models

HMM

Conditional probabilities

Structured prediction

Structured perceptron →
Classification with lots of features
over structured models!

Why are Features Good?

- Can easily try many different ideas
 - Are capital letters usually nouns?
 - Are words that end with -ed usually verbs? -ing?

Restructuring HMM With Features

Normal HMM:
$$P(X, Y) = \prod_{i=1}^l P_E(x_i | y_i) \prod_{i=1}^{l+1} P_T(y_i | y_{i-1})$$

Restructuring HMM With Features

Normal HMM:
$$P(X, Y) = \prod_{i=1}^l P_E(x_i | y_i) \prod_{i=1}^{l+1} P_T(y_i | y_{i-1})$$

Log Likelihood:
$$\log P(X, Y) = \sum_{i=1}^l \log P_E(x_i | y_i) + \sum_{i=1}^{l+1} \log P_T(y_i | y_{i-1})$$

Restructuring HMM With Features

Normal HMM:
$$P(X, Y) = \prod_{i=1}^l P_E(x_i | y_i) \prod_{i=1}^{l+1} P_T(y_i | y_{i-1})$$

Log Likelihood:
$$\log P(X, Y) = \sum_{i=1}^l \log P_E(x_i | y_i) + \sum_{i=1}^{l+1} \log P_T(y_i | y_{i-1})$$

Score
$$S(X, Y) = \sum_{i=1}^l w_{E, y_i, x_i} + \sum_{i=1}^{l+1} w_{T, y_{i-1}, y_i}$$

Restructuring HMM With Features

Normal HMM:
$$P(X, Y) = \prod_{i=1}^l P_E(x_i | y_i) \prod_{i=1}^{l+1} P_T(y_i | y_{i-1})$$

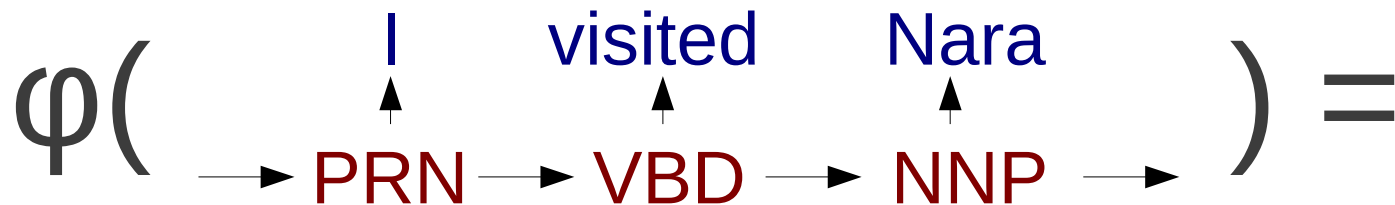
Log Likelihood:
$$\log P(X, Y) = \sum_{i=1}^l \log P_E(x_i | y_i) + \sum_{i=1}^{l+1} \log P_T(y_i | y_{i-1})$$

Score
$$S(X, Y) = \sum_{i=1}^l w_{E, y_i, x_i} + \sum_{i=1}^{l+1} w_{E, y_{i-1}, y_i}$$

When:
$$w_{E, y_i, x_i} = \log P_E(x_i | y_i) \quad w_{T, y_{i-1}, y_i} = \log P_T(y_i | y_{i-1})$$

$$\log P(X, Y) = S(X, Y)$$

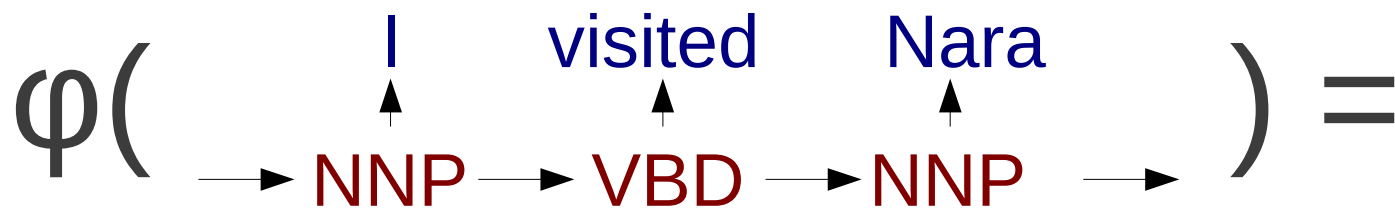
Example



$$\varphi_{T,<S>,\text{PRN}}(X, Y_1) = 1 \quad \varphi_{T,\text{PRN},\text{VBD}}(X, Y_1) = 1 \quad \varphi_{T,\text{VBD},\text{NNP}}(X, Y_1) = 1 \quad \varphi_{T,\text{NNP},</S>}(X, Y_1) = 1$$

$$\varphi_{E,\text{PRN},\text{"I"}}(X, Y_1) = 1 \quad \varphi_{E,\text{VBD},\text{"visited"}}(X, Y_1) = 1 \quad \varphi_{E,\text{NNP},\text{"Nara"}}(X, Y_1) = 1$$

$$\varphi_{\text{CAPS},\text{PRN}}(X, Y_1) = 1 \quad \varphi_{\text{CAPS},\text{NNP}}(X, Y_1) = 1 \quad \varphi_{\text{SUF},\text{VBD},\text{"...ed"}}(X, Y_1) = 1$$



$$\varphi_{T,<S>,\text{NNP}}(X, Y_1) = 1 \quad \varphi_{T,\text{NNP},\text{VBD}}(X, Y_1) = 1 \quad \varphi_{T,\text{VBD},\text{NNP}}(X, Y_1) = 1 \quad \varphi_{T,\text{NNP},</S>}(X, Y_1) = 1$$

$$\varphi_{E,\text{NNP},\text{"I"}}(X, Y_1) = 1 \quad \varphi_{E,\text{VBD},\text{"visited"}}(X, Y_1) = 1 \quad \varphi_{E,\text{NNP},\text{"Nara"}}(X, Y_1) = 1$$

$$\varphi_{\text{CAPS},\text{NNP}}(X, Y_1) = 2 \quad \varphi_{\text{SUF},\text{VBD},\text{"...ed"}}(X, Y_1) = 1$$

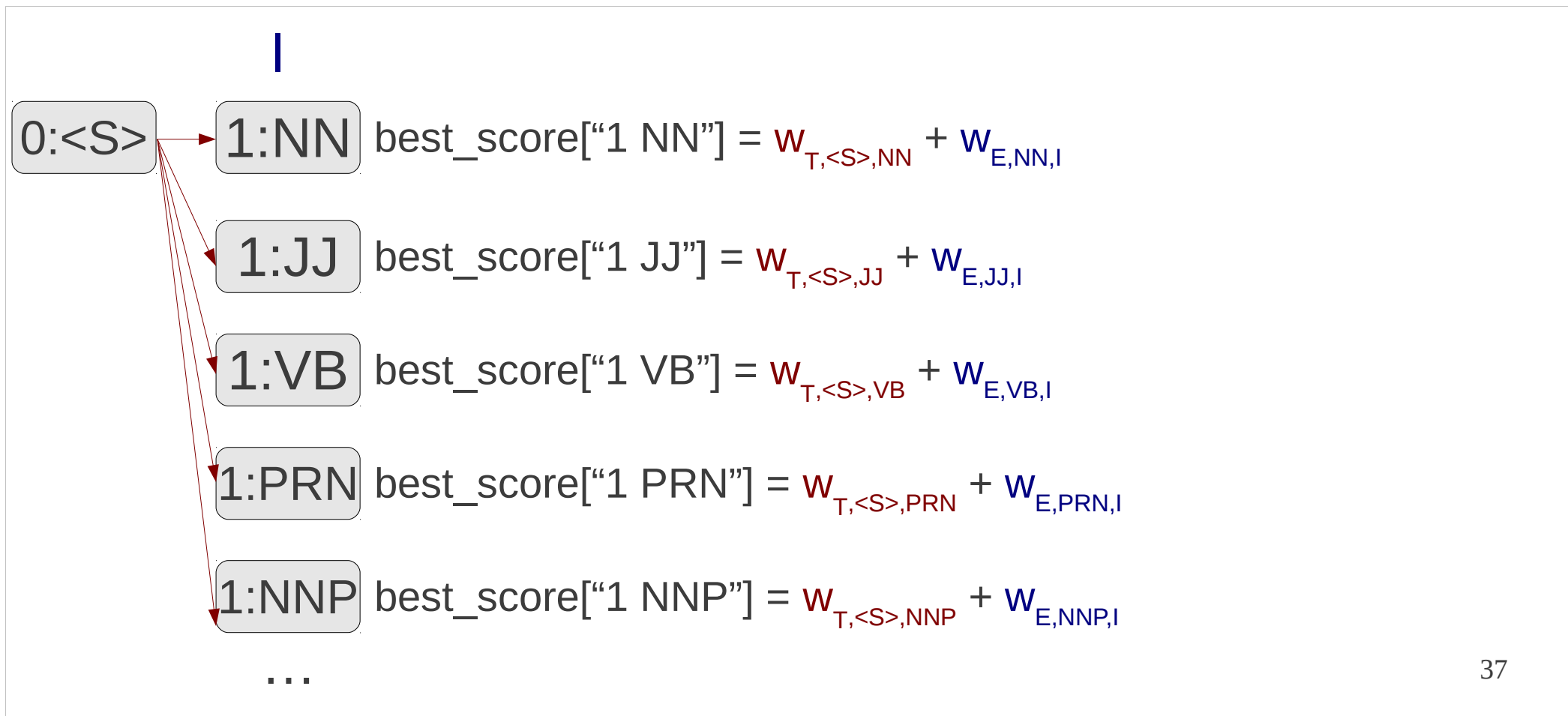
Finding the Best Solution

- We must find the POS sequence that satisfies:

$$\hat{Y} = \operatorname{argmax}_Y \sum_i w_i \varphi_i(X, Y)$$

HMM Viterbi with Features

- Same as probabilities, use feature weights



HMM Viterbi with Features

- Can add additional features



Learning In the Structured Perceptron

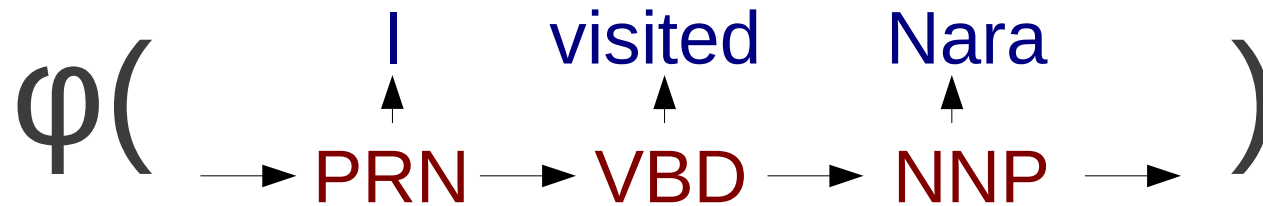
- Remember the perceptron algorithm
- If there is a mistake:

$$w \leftarrow w + y \varphi(x)$$

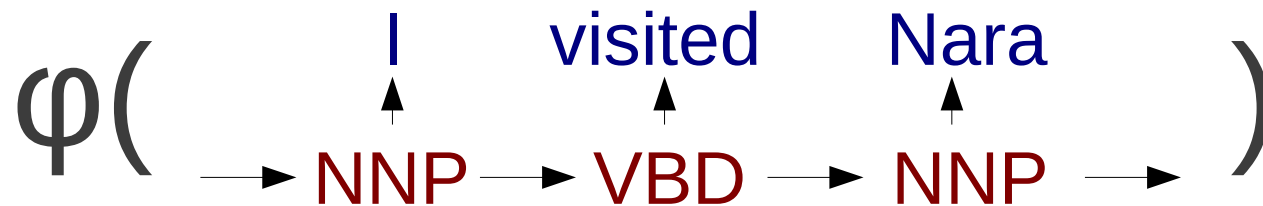
- Update weights to:
 - increase score of positive examples
 - decrease score of negative examples
- What is positive/negative in structured perceptron?

Learning in the Structured Perceptron

- Positive example, correct feature vector:

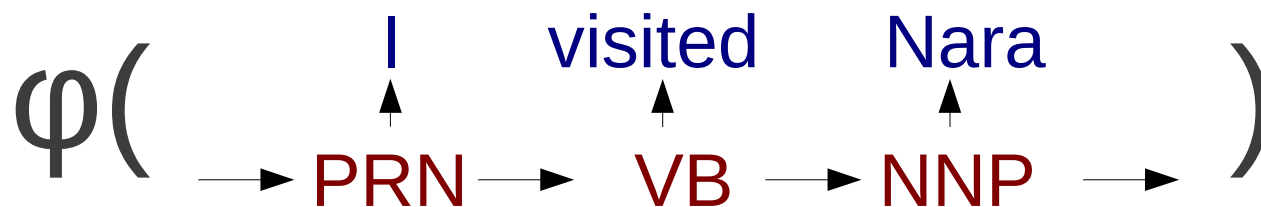
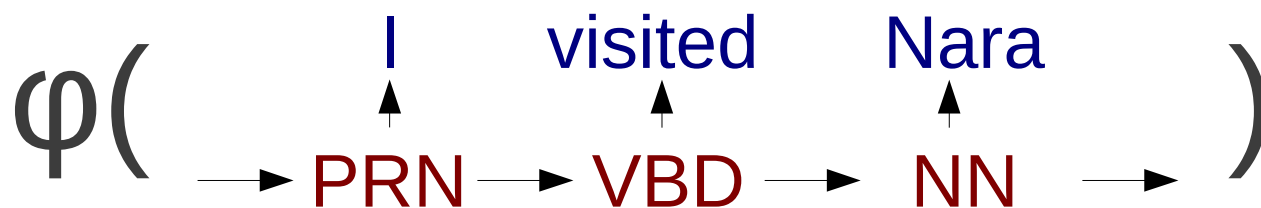
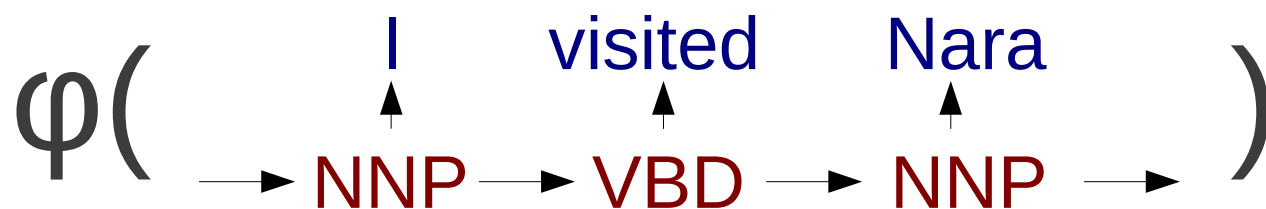


- Negative example, incorrect feature vector:



Choosing an Incorrect Feature Vector

- There are too many incorrect feature vectors!



- Which do we use?

Choosing an Incorrect Feature Vector

- Answer: We update using the incorrect answer with the highest score:

$$\hat{Y} = \operatorname{argmax}_Y \sum_i w_i \varphi_i(X, Y)$$

- Our update rule becomes:

$$\mathbf{w} \leftarrow \mathbf{w} + \varphi(X, Y') - \varphi(X, \hat{Y})$$

- (Y' is the correct answer)
- Note: If highest scoring answer is correct, no change

Example

$$\begin{array}{llll}
 \varphi_{T,<S>,PRN}(X,Y_1) = 1 & \varphi_{T,PRN,VBD}(X,Y_1) = 1 & \varphi_{T,VBD,NNP}(X,Y_1) = 1 & \varphi_{T,NNP,</S>}(X,Y_1) = 1 \\
 \varphi_{E,PRN,"I"}(X,Y_1) = 1 & \varphi_{E,VBD,"visited"}(X,Y_1) = 1 & \varphi_{E,NNP,"Nara"}(X,Y_1) = 1 & \\
 \varphi_{CAPS,PRN}(X,Y_1) = 1 & \varphi_{CAPS,NNP}(X,Y_1) = 1 & \varphi_{SUF,VBD,"...ed"}(X,Y_1) = 1 &
 \end{array}$$

–

$$\begin{array}{llll}
 \varphi_{T,<S>,NNP}(X,Y_1) = 1 & \varphi_{T,NNP,VBD}(X,Y_1) = 1 & \varphi_{T,VBD,NNP}(X,Y_1) = 1 & \varphi_{T,NNP,</S>}(X,Y_1) = 1 \\
 \varphi_{E,NNP,"I"}(X,Y_1) = 1 & \varphi_{E,VBD,"visited"}(X,Y_1) = 1 & \varphi_{E,NNP,"Nara"}(X,Y_1) = 1 & \\
 \varphi_{CAPS,NNP}(X,Y_1) = 2 & & \varphi_{SUF,VBD,"...ed"}(X,Y_1) = 1 &
 \end{array}$$

=

$$\begin{array}{llll}
 \varphi_{T,<S>,PRN}(X,Y_1) = 1 & \varphi_{T,PRN,VBD}(X,Y_1) = 1 & \varphi_{T,VBD,NNP}(X,Y_1) = 0 & \varphi_{T,NNP,</S>}(X,Y_1) = 0 \\
 \varphi_{T,<S>,NNP}(X,Y_1) = -1 & \varphi_{T,NNP,VBD}(X,Y_1) = -1 & & \\
 \varphi_{E,PRN,"I"}(X,Y_1) = 1 & \varphi_{E,VBD,"visited"}(X,Y_1) = 0 & \varphi_{E,NNP,"Nara"}(X,Y_1) = 0 & \\
 \varphi_{E,NNP,"I"}(X,Y_1) = -1 & & & \\
 \varphi_{CAPS,NNP}(X,Y_1) = -1 & \varphi_{CAPS,PRN}(X,Y_1) = 1 & \varphi_{SUF,VBD,"...ed"}(X,Y_1) = 0 &
 \end{array}$$

Structured Perceptron Algorithm

```
create map  $w$ 
for / iterations
  for each labeled pair  $X$ ,  $Y_{prime}$  in the data
     $Y_{hat} = \text{HMM\_VITERBI}(w, X)$ 
     $\phi_{prime} = \text{CREATE\_FEATURES}(X, Y_{prime})$ 
     $\phi_{hat} = \text{CREATE\_FEATURES}(X, Y_{hat})$ 
     $w += \phi_{prime} - \phi_{hat}$ 
```

Conclusion

- The structured perceptron is a **discriminative structured** prediction model
 - HMM: generative structured prediction
 - Perceptron: discriminative binary prediction
- It can be used for **many** problems
 - Prediction of

Thank You!