

Learning for Structured Prediction

Linear Methods For Sequence Labeling:  
Hidden Markov Models vs Structured Perceptron

Ivan Titov

$\mathbf{x}$  is an input (e.g., sentence),  $\mathbf{y}$  is an output (syntactic tree)

# Last Time: Structured Prediction

## 1. Selecting feature representation $\varphi(\mathbf{x}, \mathbf{y})$

- ▶ It should be sufficient to discriminate correct structure from incorrect ones
- ▶ It should be possible to decode with it (see (3))

## 2. Learning

- ▶ Which error function to optimize on the training set, for example

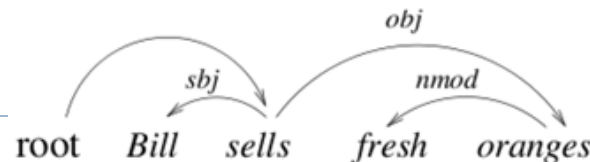
$$w \cdot \varphi(\mathbf{x}, \mathbf{y}^*) - \max_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x}), \mathbf{y}' \neq \mathbf{y}^*} w \cdot \varphi(\mathbf{x}, \mathbf{y}') > \gamma$$

- ▶ How to make it efficient (see (3))

## 3. Decoding: $\mathbf{y} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} w \cdot \varphi(\mathbf{x}, \mathbf{y}')$

- ▶ Dynamic programming for simpler representations  $\varphi$  ?
- ▶ Approximate search for more powerful ones?

We illustrated all these challenges on the example of dependency parsing



# Outline

---

- ▶ Sequence labeling / segmentation problems: settings and example problems:
  - ▶ Part-of-speech tagging, named entity recognition, gesture recognition
- ▶ Hidden Markov Model
  - ▶ Standard definition + maximum likelihood estimation
  - ▶ General views: as a representative of linear models
- ▶ Perceptron and Structured Perceptron
  - ▶ algorithms / motivations
- ▶ Decoding with the Linear Model
- ▶ Discussion: Discriminative vs. Generative

# Sequence Labeling Problems

---

- ▶ Definition:

- ▶ Input: sequences of variable length  $\mathbf{x} = (x_1, x_2, \dots, x_{|\mathbf{x}|}), x_i \in X$

- ▶ Output: every position is labeled  $\mathbf{y} = (y_1, y_2, \dots, y_{|\mathbf{x}|}), y_i \in Y$

- ▶ Examples:

- ▶ Part-of-speech tagging

|                |             |                |          |            |            |          |
|----------------|-------------|----------------|----------|------------|------------|----------|
| $\mathbf{x} =$ | <b>John</b> | <b>carried</b> | <b>a</b> | <b>tin</b> | <b>can</b> | <b>.</b> |
| $\mathbf{y} =$ | NP          | VBD            | DT       | NN         | NN         | .        |

- ▶ Named-entity recognition, shallow parsing (“chunking”), gesture recognition from video-streams, ...

# Part-of-speech tagging

|       |             |                |          |            |            |   |
|-------|-------------|----------------|----------|------------|------------|---|
| $x =$ | <b>John</b> | <b>carried</b> | <b>a</b> | <b>tin</b> | <b>can</b> | . |
| $y =$ | NNP         | VBD            | DT       | NN         | NN or MD?  | . |

## ► Labels:

- NNP – proper singular noun;
- VBD - verb, past tense
- DT - determiner
- NN – singular noun
- MD - modal

## ► Consider

|       |            |            |              |                  |     |
|-------|------------|------------|--------------|------------------|-----|
| $x =$ | <b>Tin</b> | <b>can</b> | <b>cause</b> | <b>poisoning</b> | ... |
| $y =$ | NN         | MD         | VB           | NN               | ... |

In fact, even knowing that the previous word is a noun is not enough. The most word you will make a mistake here

One need to model interactions between labels to successfully resolve ambiguities, so this should be tackled as a structured prediction problem

# Named Entity Recognition

[**ORG** Chelsea], despite their name, are not based in [**LOC** Chelsea], but in neighbouring [**LOC** Fulham] .

- ▶ Not as trivial as it may seem, consider:

Chelsea can be a person too!

- ▶ [PERS Bill Clinton] will not embarrass [PERS **Chelsea**] at her wedding
- ▶ **Tiger** failed to make a birdie in the South Course ...

- ▶ Encoding example (BIO-encoding)

Is it an animal or a person?

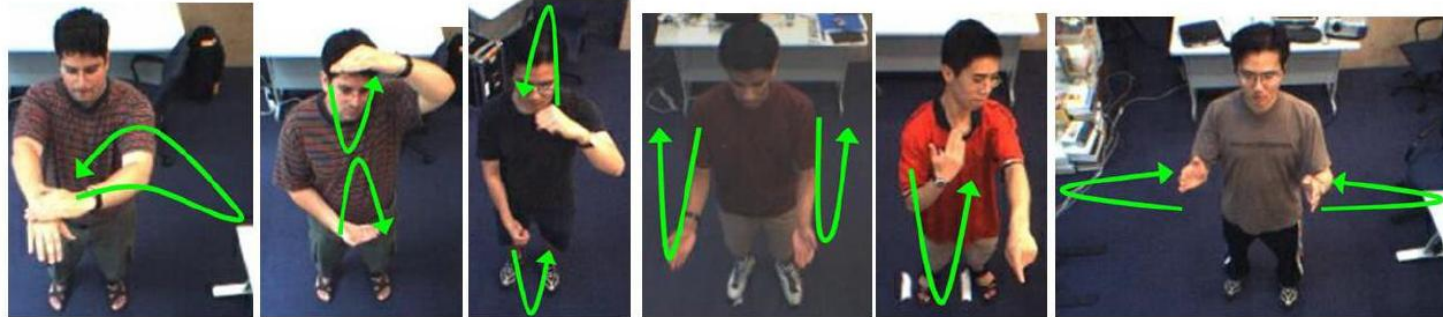
|       |             |                |                    |                |           |            |                |           |              |               |
|-------|-------------|----------------|--------------------|----------------|-----------|------------|----------------|-----------|--------------|---------------|
| $x =$ | <b>Bill</b> | <b>Clinton</b> | <b>embarrassed</b> | <b>Chelsea</b> | <b>at</b> | <b>her</b> | <b>wedding</b> | <b>at</b> | <b>Astor</b> | <b>Courts</b> |
| $y =$ | B-PERS      | I-PERS         | O                  | B-PERS         | O         | O          | O              | O         | B-LOC        | I-LOC         |

# Vision: Gesture Recognition

- ▶ Given a sequence of frames in a video annotate each frame with a gesture type:



- ▶ Types of gestures:



Flip back

Shrink  
vertically

Expand  
vertically

Double  
back

Point  
and back

Expand horizontally

- ▶ It is hard to predict gestures from each frame in isolation, you need to **exploit relations** between frames and gesture types

# Outline

---

- ▶ Sequence labeling / segmentation problems: settings and example problems:
  - ▶ Part-of-speech tagging, named entity recognition, gesture recognition
- ▶ Hidden Markov Model
  - ▶ Standard definition + maximum likelihood estimation
  - ▶ General views: as a representative of linear models
- ▶ Perceptron and Structured Perceptron
  - ▶ algorithms / motivations
- ▶ Decoding with the Linear Model
- ▶ Discussion: Discriminative vs. Generative



# Hidden Markov Models

---

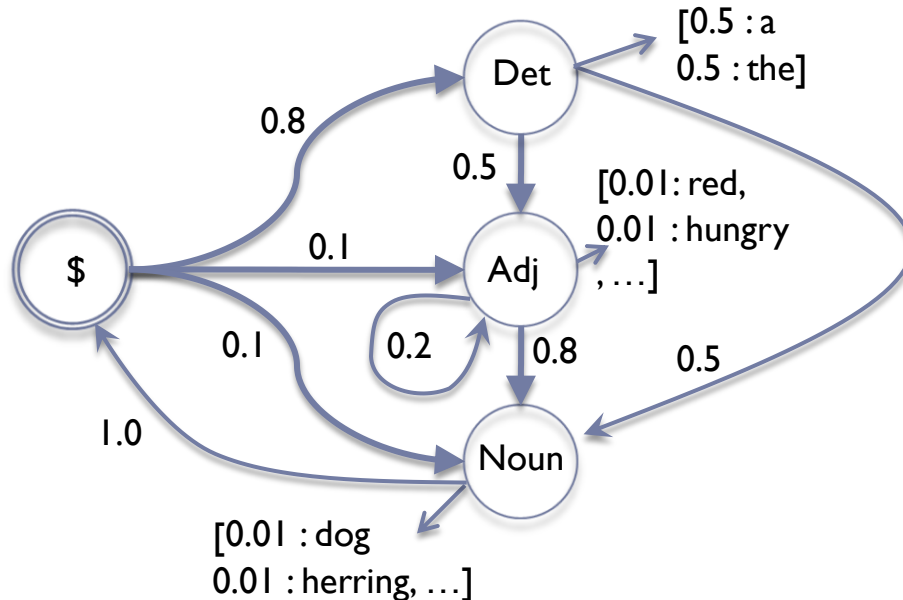
- ▶ We will consider the part-of-speech (POS) tagging example

|             |                |          |            |            |          |
|-------------|----------------|----------|------------|------------|----------|
| <b>John</b> | <b>carried</b> | <b>a</b> | <b>tin</b> | <b>can</b> | <b>.</b> |
| NP          | VBD            | DT       | NN         | NN         | .        |

- ▶ A “generative” model, i.e.:
  - ▶ Model: Introduce a parameterized model of how both words and tags are generated  $P(x, y|\theta)$
  - ▶ Learning: use a labeled training set to estimate the most likely parameters of the model  $\hat{\theta}$
  - ▶ Decoding:  $y = \operatorname{argmax}_{y'} P(x, y'|\hat{\theta})$

# Hidden Markov Models

A simplistic state diagram for noun phrases:       $N$  – tags,    $M$  – vocabulary size



Example:

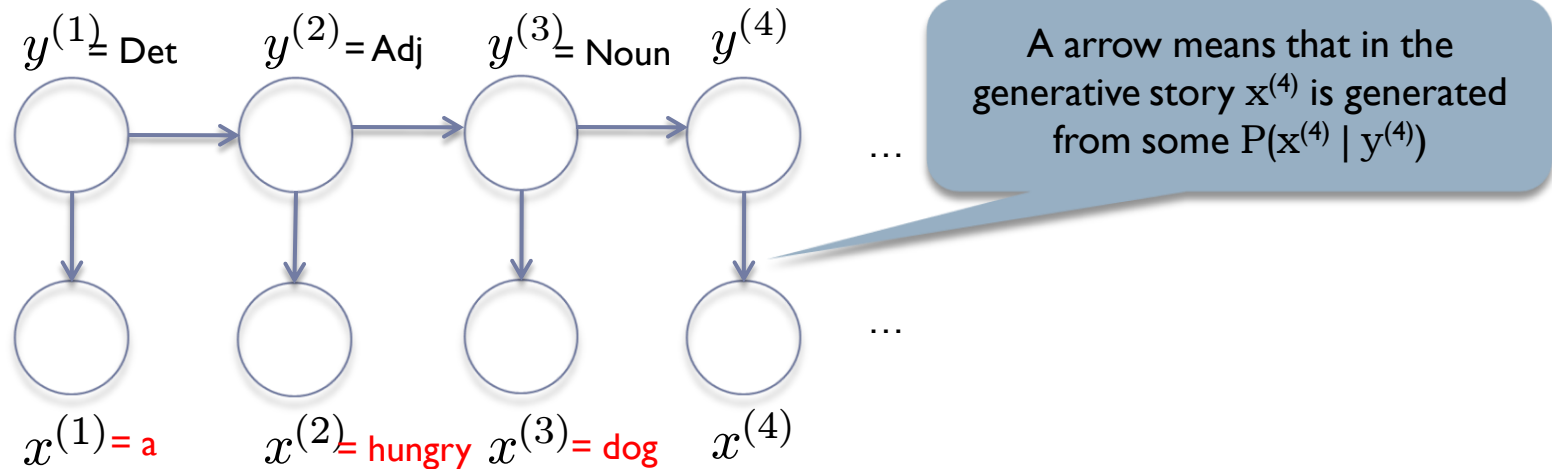
a      hungry      dog

- ▶ States correspond to POS tags,
- ▶ Words are emitted independently from each POS tag
- ▶ Parameters (to be estimated from the training set):
  - ▶ Transition probabilities  $P(y^{(t)}|y^{(t-1)})$  :  $[N \times N]$  matrix
  - ▶ Emission probabilities  $P(x^{(t)}|y^{(t)})$  :  $[N \times M]$  matrix

Stationarity assumption: this probability does not depend on the position in the sequence  $t$

# Hidden Markov Models

Representation as an instantiation of a *graphical model*:  $N$  – tags,  $M$  – vocabulary size



- ▶ States correspond to POS tags,
- ▶ Words are emitted independently from each POS tag
- ▶ Parameters (to be estimated from the training set):
  - ▶ Transition probabilities  $P(y^{(t)} | y^{(t-1)})$  :  $[N \times N]$  matrix
  - ▶ Emission probabilities  $P(x^{(t)} | y^{(t)})$  :  $[N \times M]$  matrix

Stationarity assumption: this probability does not depend on the position in the sequence  $t$

# Hidden Markov Models: Estimation

---

- ▶  $N$  – the number tags,  $M$  – vocabulary size
- ▶ Parameters (to be estimated from the training set):
  - ▶ Transition probabilities  $a_{ji} = P(y^{(t)} = i | y^{(t-1)} = j)$ ,  $A$  -  $[N \times N]$  matrix
  - ▶ Emission probabilities  $b_{ik} = P(x^{(t)} = k | y^{(t)} = i)$ ,  $B$  -  $[N \times M]$  matrix
- ▶ Training corpus:
  - ▶  $x^{(1)} = (\text{In, an, Oct., 19, review, of, ...}), y^{(1)} = (\text{IN, DT, NNP, CD, NN, IN, ...})$
  - ▶  $x^{(2)} = (\text{Ms., Haag, plays, Elianti,.}), y^{(2)} = (\text{NNP, NNP, VBZ, NNP,.})$
  - ▶ ...
  - ▶  $x^{(L)} = (\text{The, company, said,...}), y^{(L)} = (\text{DT, NN, VBD, NNP,.})$
- ▶ How to estimate the parameters using maximum likelihood estimation?
  - ▶ You probably can guess what these estimation should be?

# Hidden Markov Models: Estimation

- Parameters (to be estimated from the training set):
  - Transition probabilities  $a_{ji} = P(y^{(t)} = i | y^{(t-1)} = j)$ ,  $A$  -  $[N \times N]$  matrix
  - Emission probabilities  $b_{ik} = P(x^{(t)} = k | y^{(t)} = i)$ ,  $B$  -  $[N \times M]$  matrix
- Training corpus:  $(\mathbf{x}^{(l)}, \mathbf{y}^{(l)})$ ,  $l = 1, \dots, L$
- Write down the probability of the corpus according to the HMM:

$$P(\{\mathbf{x}^{(l)}, \mathbf{y}^{(l)}\}_{l=1}^L) = \prod_{l=1}^L P(\mathbf{x}^{(l)}, \mathbf{y}^{(l)}) =$$

$$= \prod_{l=1}^L a_{\$,y_1^{(l)}} \left( \prod_{t=1}^{|\mathbf{x}^{(l)}|-1} b_{y_t^{(l)}, x_t^{(l)}} a_{y_t^{(l)}, y_{t+1}^{(l)}} \right) b_{y_{|\mathbf{x}^{(l)}|}^{(l)}, x_{|\mathbf{x}^{(l)}|}^{(l)}} a_{y_{|\mathbf{x}^{(l)}|}^{(l)}, \$} =$$

Select tag for  
the first word

Draw a word  
from this state

Select the next  
state

Draw last  
word

Transit into the  
\$ state

$$= \prod_{i,j=1}^N a_{i,j}^{C_T(i,j)} \prod_{i=1}^N \prod_{k=1}^M b_{i,k}^{C_E(i,k)}$$

$C_T(i,j)$  is #times tag  $i$  is followed by tag  $j$ .  
Here we assume that  $\$$  is a special tag which  
precedes and succeeds every sentence

$C_E(i,k)$  is #times  
word  $k$  is  
emitted by tag  $i$

# Hidden Markov Models: Estimation

► Maximize:  $P(\{\mathbf{x}^{(l)}, \mathbf{y}^{(l)}\}_{l=1}^L) =$   
 $= \prod_{i,j=1}^N a_{i,j}^{C_T(i,j)} \prod_{i=1}^N \prod_{k=1}^M b_{i,k}^{C_E(i,k)}$

$C_T(i,j)$  is #times tag  $i$  is followed by tag  $j$ .

$C_E(i,k)$  is #times word  $k$  is emitted by tag  $i$

► Equivalently maximize the logarithm of this:  $\log(P(\{\mathbf{x}^{(l)}, \mathbf{y}^{(l)}\}_{l=1}^L)) =$   
 $= \sum_{i=1}^N \left( \sum_{j=1}^N C_T(i,j) \log a_{i,j} + \sum_{k=1}^M C_E(i,k) \log b_{i,k} \right)$

subject to probabilistic constraints:  $\sum_{j=1}^N a_{i,j} = 1, \quad \sum_{k=1}^M b_{i,k} = 1, \quad i = 1, \dots, N$

► Or, we can decompose it into  $2N$  optimization tasks:

For transitions

$i = 1, \dots, N:$

$$\max_{a_{i,1}, \dots, a_{i,N}} \sum_{j=1}^N C_T(i,j) \log a_{i,j}$$

$$\text{s.t. } \sum_{j=1}^N a_{i,j} = 1$$

For emissions

$i = 1, \dots, N:$

$$\max_{b_{i,1}, \dots, b_{i,M}} C_E(i,k) \log b_{i,k}$$

$$\text{s.t. } \sum_{k=1}^M b_{i,k} = 1$$

# Hidden Markov Models: Estimation

- ▶ For transitions (some  $i$ )

$$\max_{a_{i,1}, \dots, a_{i,N}} \sum_{j=1}^N C_T(i, j) \log a_{i,j}$$

$$\text{s.t. } 1 - \sum_{j=1}^N a_{i,j} = 0$$

- ▶ Constrained optimization task, Lagrangian:

$$L(a_{i,1}, \dots, a_{i,N}, \lambda) = \sum_{j=1}^N C_T(i, j) \log a_{i,j} + \lambda \times (1 - \sum_{j=1}^N a_{i,j})$$

- ▶ Find critical points of Lagrangian by solving the system of equations:

$$\frac{\partial L}{\partial \lambda} = 1 - \sum_{j=1}^N a_{i,j} = 0$$

$$\frac{\partial L}{\partial a_{i,j}} = \frac{C_T(i,j)}{a_{i,j}} - \lambda = 0 \implies a_{i,j} = \frac{C_T(i,j)}{\lambda}$$

$$P(y^t = j | y^{t-1} = i) = a_{i,j} = \frac{C_T(i,j)}{\sum_{j'} C_T(i,j')}$$

- ▶ Similarly, for emissions:

$$P(x^t = k | y^t = i) = b_{i,k} = \frac{C_E(i,k)}{\sum_{k'} C_E(i,k')}$$

The maximum likelihood solution is just normalized counts of events.  
Always like this for generative models if all the labels are visible in training

I ignore “smoothing” to process rare or unseen word tag combinations...  
Outside score of the seminar

# HMMs as linear models

|             |                |          |            |            |          |
|-------------|----------------|----------|------------|------------|----------|
| <b>John</b> | <b>carried</b> | <b>a</b> | <b>tin</b> | <b>can</b> | <b>.</b> |
| ?           | ?              | ?        | ?          | ?          | .        |

- ▶ Decoding:  $\mathbf{y} = \operatorname{argmax}_{\mathbf{y}'} P(\mathbf{x}, \mathbf{y}' | A, B) = \operatorname{argmax}_{\mathbf{y}'} \log P(\mathbf{x}, \mathbf{y}' | A, B)$
- ▶ We will talk about the decoding algorithm slightly later, let us generalize Hidden Markov Model:

$$\begin{aligned} \log P(\mathbf{x}, \mathbf{y}' | A, B) &= \sum_{l=1}^{|\mathbf{x}|+1} \log b_{y'_l, x_l} + \log a_{y'_l, y'_{l+1}} \\ &= \sum_{i=1}^N \sum_{j=1}^N C_T(\mathbf{y}', i, j) \times \log a_{i,j} + \sum_{i=1}^N \sum_{k=1}^M C_E(\mathbf{x}, \mathbf{y}', i, k) \times \log b_{i,k} \end{aligned}$$

The number of times tag  $i$  is followed by tag  $j$  in the candidate  $\mathbf{y}'$

The number of times tag  $i$  corresponds to word  $k$  in  $(\mathbf{x}, \mathbf{y}')$

- ▶ But this is just a linear model!!



# Scoring: example

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

|             |                |          |            |            |          |
|-------------|----------------|----------|------------|------------|----------|
| <b>John</b> | <b>carried</b> | <b>a</b> | <b>tin</b> | <b>can</b> | <b>.</b> |
| NP          | VBD            | DT       | NN         | NN         | .        |

$\varphi(\mathbf{x}, \mathbf{y}') = ($ 

|   |   |     |   |   |   |     |
|---|---|-----|---|---|---|-----|
| 1 | 0 | ... | 1 | 1 | 0 | ... |
|---|---|-----|---|---|---|-----|

 $)$

Unary features  
 $C_E(\mathbf{x}, \mathbf{y}', i, k)$

NP: John   NP: Mary   ...   NP-VBD   NN-.   MD-.   ...

Edge features  
 $C_T(\mathbf{y}', i, j)$

$w_{ML} = ($ 

|                     |                     |     |                    |                  |                  |     |
|---------------------|---------------------|-----|--------------------|------------------|------------------|-----|
| $\log b_{NP, John}$ | $\log b_{NP, Mary}$ | ... | $\log a_{NN, VBD}$ | $\log a_{NN, .}$ | $\log a_{MD, .}$ | ... |
|---------------------|---------------------|-----|--------------------|------------------|------------------|-----|

 $)$

- ▶ Their inner product is exactly  $\log P(\mathbf{x}, \mathbf{y}' | A, B)$

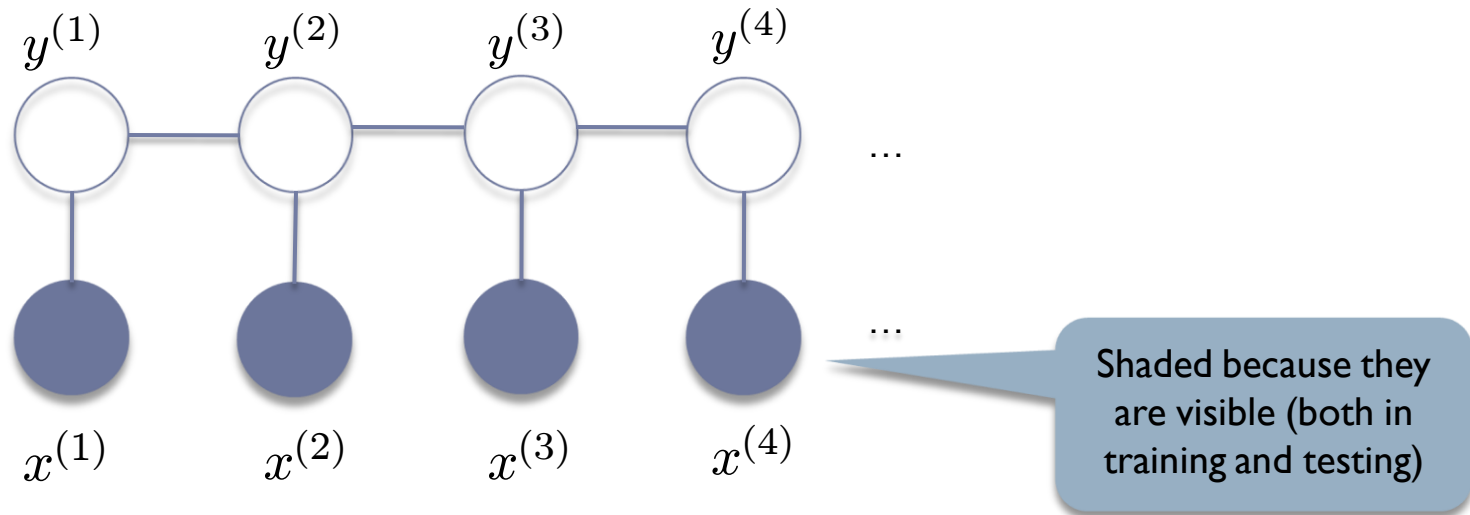
$$w_{ML} \cdot \varphi(\mathbf{x}, \mathbf{y}') = \sum_{i=1}^N \sum_{j=1}^N C_T(\mathbf{y}', i, j) \times \log a_{i,j} + \sum_{i=1}^N \sum_{k=1}^M C_E(\mathbf{x}, \mathbf{y}', i, k) \times \log b_{i,k}$$

- ▶ But may be there other (and better?) ways to estimate  $w$ , especially when we know that HMM is not a faithful model of reality?
- ▶ **It is not only a theoretical question!** (we'll talk about that in a moment)



# Feature view

Basically, we define features which correspond to edges in the graph:

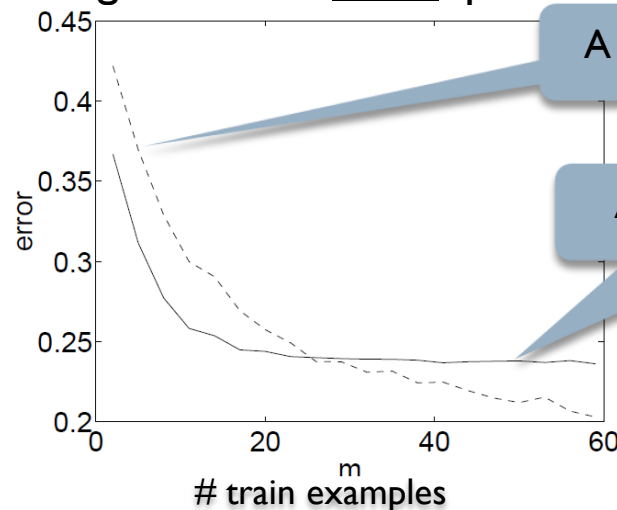


# Generative modeling

Real case: HMM is a coarse approximation of reality

- ▶ For a very large dataset (asymptotic analysis):
  - ▶ If data is generated from some “true” HMM, then (if the training set is sufficiently large), we are guaranteed to have an optimal tagger
  - ▶ Otherwise, (generally) HMM will not correspond to an optimal linear classifier
  - ▶ Discriminative methods which minimize the error more directly **are guaranteed** (under some fairly general conditions) to converge to an optimal linear classifier
- ▶ For smaller training sets
  - ▶ Generative classifiers converge faster to their optimal error [Ng & Jordan, NIPS 01]

Errors on a regression dataset (predict housing prices in Boston area):



A discriminative classifier

A generative model

# Outline

---

- ▶ Sequence labeling / segmentation problems: settings and example problems:
  - ▶ Part-of-speech tagging, named entity recognition, gesture recognition
- ▶ Hidden Markov Model
  - ▶ Standard definition + maximum likelihood estimation
  - ▶ General views: as a representative of linear models
- ▶ Perceptron and Structured Perceptron
  - ▶ algorithms / motivations
- ▶ Decoding with the Linear Model
- ▶ Discussion: Discriminative vs. Generative

# Perceptron

- ▶ Let us start with a binary classification problem  $y \in \{+1, -1\}$ 
  - ▶ For binary classification the prediction rule is:  $y = \text{sign}(\mathbf{w} \cdot \varphi(x))$
- ▶ Perceptron algorithm, given a training set  $\{x^{(l)}, y^{(l)}\}_{l=1}^L$

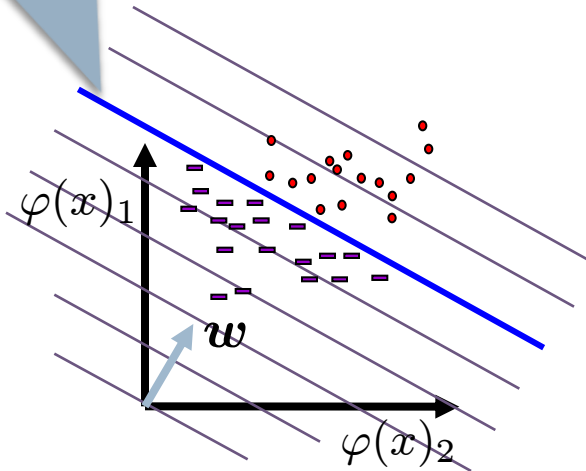
*break ties (0) in some deterministic way*

```
 $\mathbf{w} = \mathbf{0}$  // initialize
do
  err = 0
  for  $l = 1 \dots L$  // over the training examples
    if (  $y^{(l)} (\mathbf{w} \cdot \varphi(x^{(l)})) < 0$  ) // if mistake
       $\mathbf{w} += \eta y^{(l)} \varphi(x^{(l)})$  // update,  $\eta > 0$ 
      err ++ // # errors
    endif
  endfor
while ( err > 0 ) // repeat until no errors
return  $\mathbf{w}$ 
```

# Linear classification

- ▶ Linear separable case, “a perfect” classifier:

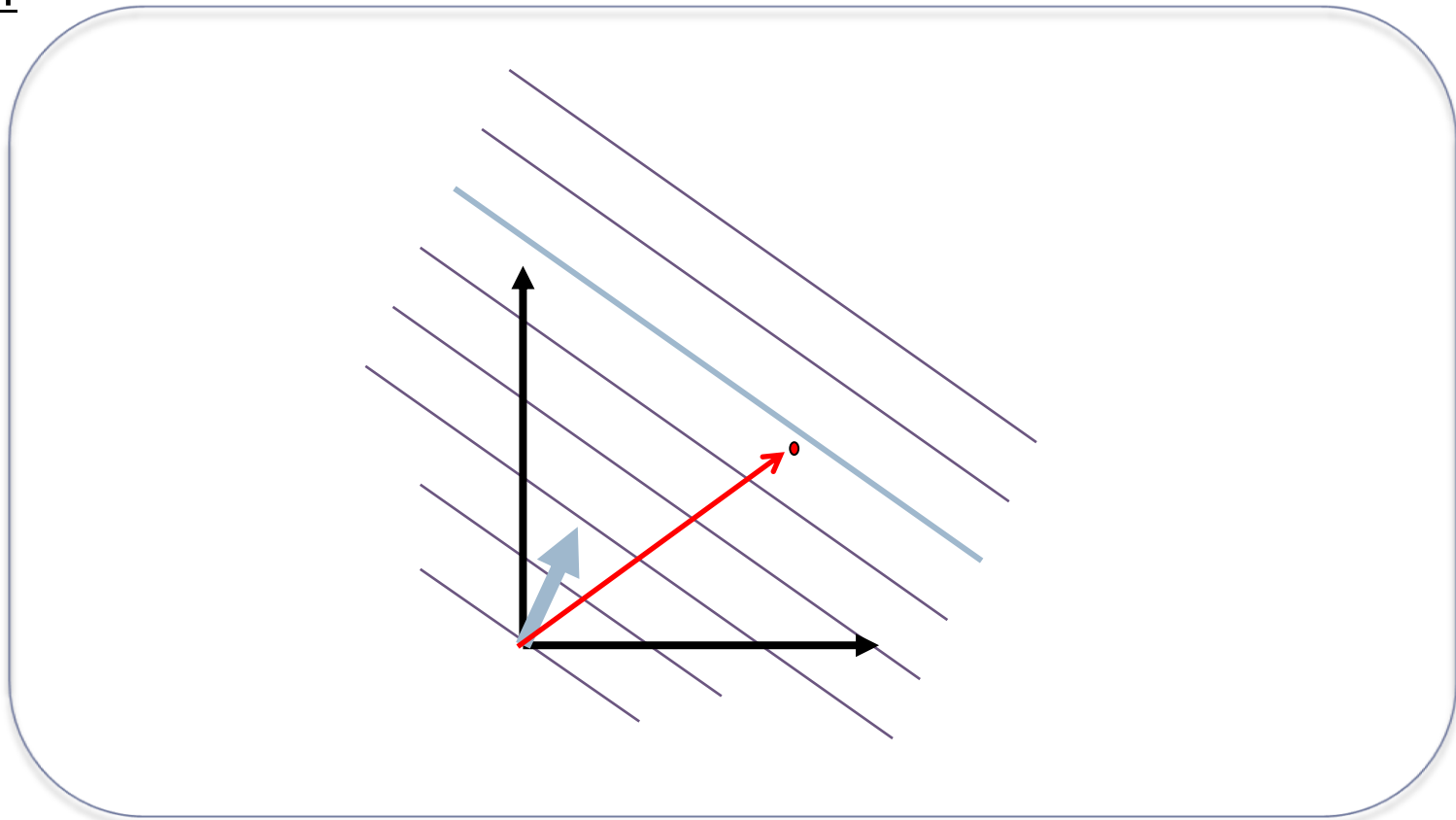
$$(\mathbf{w} \cdot \varphi(x) + b) = 0$$



- ▶ Linear functions are often written as:  $y = \text{sign}(\mathbf{w} \cdot \varphi(x) + b)$ , but we can assume that  $\varphi(x)_0 = 1$  for any  $x$

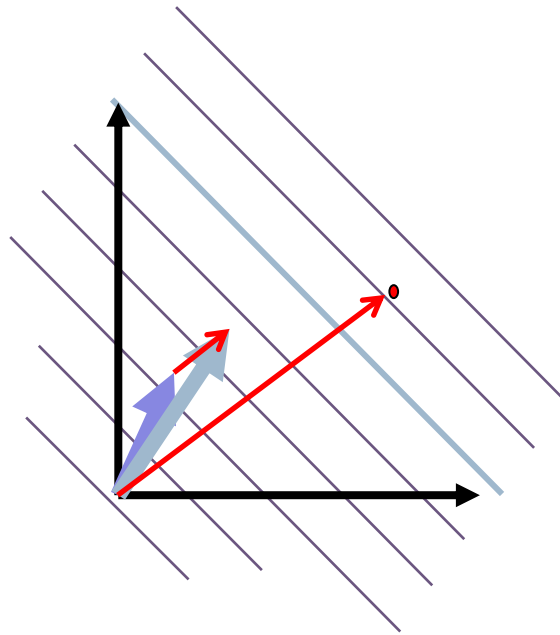
# Perceptron: geometric interpretation

```
if (  $y^{(l)} (w \cdot \varphi(x^{(l)})) < 0$  ) // if mistake  
     $w \mathrel{+}= \eta y^{(l)} \varphi(x^{(l)})$  // update  
endif
```



# Perceptron: geometric interpretation

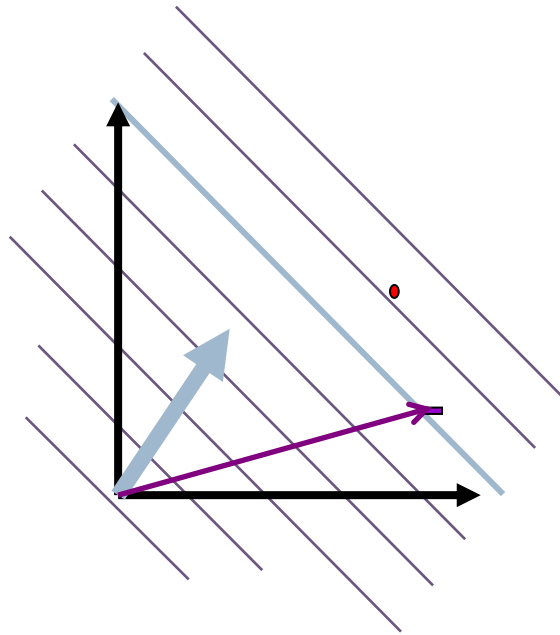
```
if (  $y^{(l)} (\mathbf{w} \cdot \varphi(x^{(l)})) < 0$  ) // if mistake  
     $\mathbf{w} += \eta y^{(l)} \varphi(x^{(l)})$  // update  
endif
```





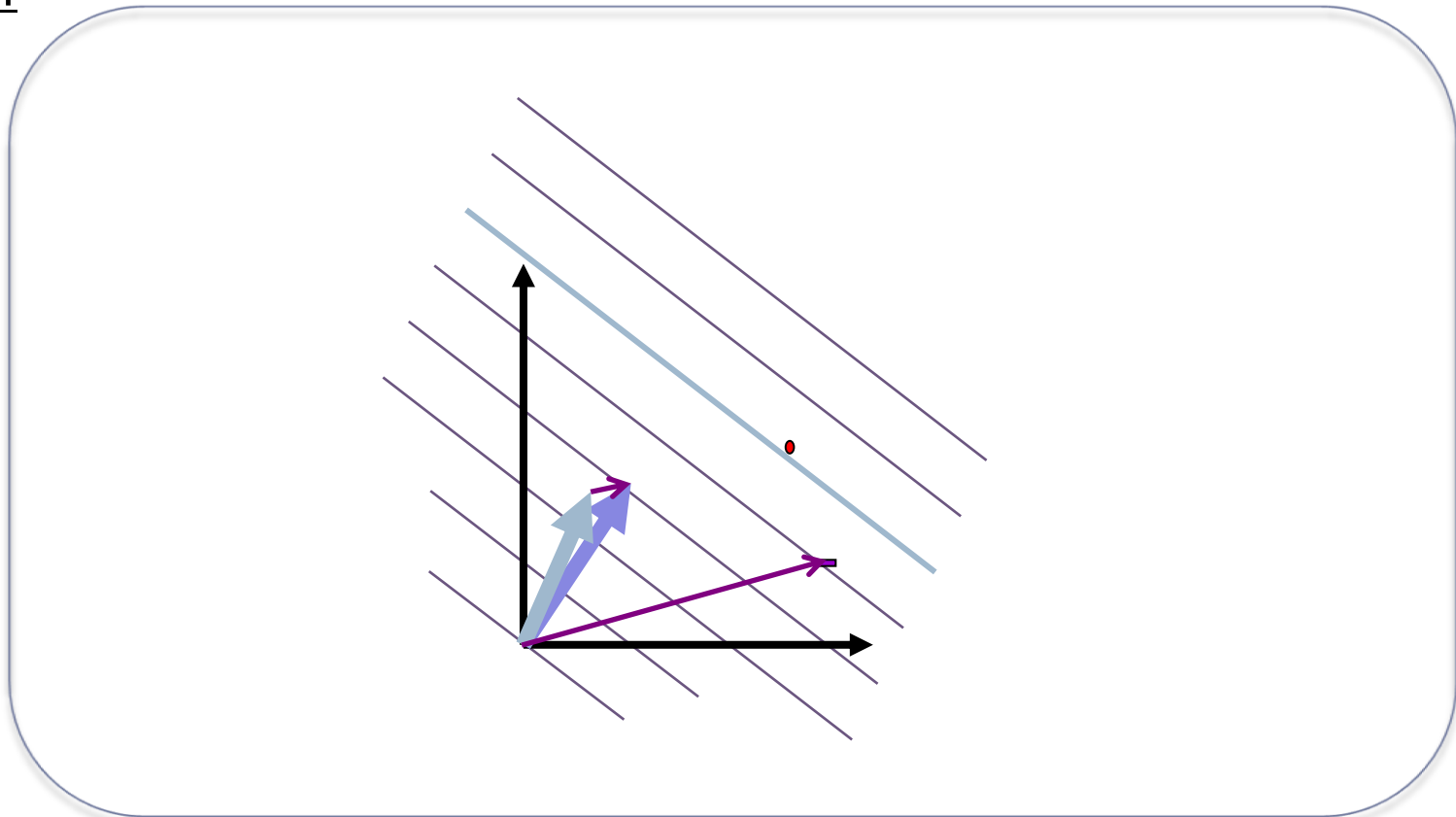
# Perceptron: geometric interpretation

```
if (  $y^{(l)} (w \cdot \varphi(x^{(l)})) < 0$  ) // if mistake  
     $w \mathrel{+}= \eta y^{(l)} \varphi(x^{(l)})$  // update  
endif
```



# Perceptron: geometric interpretation

```
if (  $y^{(l)} (w \cdot \varphi(x^{(l)})) < 0$  ) // if mistake  
     $w \mathrel{+}= \eta y^{(l)} \varphi(x^{(l)})$  // update  
endif
```



# Perceptron: algebraic interpretation

```
if (  $y^{(l)}(\mathbf{w} \cdot \varphi(x^{(l)})) < 0$  )    // if mistake  
     $\mathbf{w} += \eta y^{(l)} \varphi(x^{(l)})$     // update  
endif
```

- ▶ We want after the update to increase  $y^{(l)}(\mathbf{w} \cdot \varphi(x^{(l)}))$ 
  - ▶ If the increase is large enough than there will be no misclassification
- ▶ Let's see that's what happens after the update

$$\begin{aligned} & y^{(l)}((\mathbf{w} + \eta y^{(l)} \varphi(x^{(l)})) \cdot \varphi(x^{(l)})) \\ &= y^{(l)}(\mathbf{w} \cdot \varphi(x^{(l)})) + \eta (y^{(l)})^2 (\varphi(x^{(l)}) \cdot \varphi(x^{(l)})) \end{aligned}$$

$$(y^{(l)})^2 = 1$$

*squared norm > 0*

- ▶ So, the perceptron update moves the decision hyperplane towards misclassified  $\varphi(x^{(l)})$

# Perceptron

---

- ▶ The perceptron algorithm, obviously, can only converge if the training set is linearly separable
- ▶ It is guaranteed to converge in a finite number of iterations, dependent on how well two classes are separated (Novikoff, 1962)

# Averaged Perceptron

- ▶ A small modification

```
 $w = \mathbf{0}, w_{\Sigma} = \mathbf{0}$  // initialize  
for  $k = 1 \dots K$  // for a number of iterations  
    for  $l = 1 \dots L$  // over the training examples  
        if (  $y^{(l)}(w \cdot \varphi(x^{(l)})) < 0$  ) // if mistake  
             $w \mathrel{+}= \eta y^{(l)} \varphi(x^{(l)})$  // update,  $\eta > 0$   
        endif  
         $w_{\Sigma} \mathrel{+}= w$  // sum of  $w$  over the course of training  
    endfor  
endfor  
return  $\frac{1}{KL} w_{\Sigma}$ 
```

Do not run until convergence

Note: it is after endif

More stable in training: a vector  $w$  which survived more iterations without updates is more similar to the resulting vector  $\frac{1}{KL} w_{\Sigma}$ , as it was added a larger number of times

# Structured Perceptron

- ▶ Let us start with structured problem:  $\mathbf{y} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} \mathbf{w} \cdot \varphi(\mathbf{x}, \mathbf{y}')$
- ▶ Perceptron algorithm, given a training set  $\{\mathbf{x}^{(l)}, \mathbf{y}^{(l)}\}_{l=1}^L$

```
 $\mathbf{w} = 0$  // initialize
do
    err = 0
    for  $l = 1 \dots L$  // over the training examples
         $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x}^{(l)})} \mathbf{w} \cdot \varphi(\mathbf{x}^{(l)}, \mathbf{y}')$  // model prediction
        if (  $\mathbf{w} \cdot \varphi(\mathbf{x}^{(l)}, \hat{\mathbf{y}}) > \mathbf{w} \cdot \varphi(\mathbf{x}^{(l)}, \mathbf{y}^{(l)})$  ) // if mistake
             $\mathbf{w} += \eta (\varphi(\mathbf{x}^{(l)}, \mathbf{y}^{(l)}) - \varphi(\mathbf{x}^{(l)}, \hat{\mathbf{y}}))$  // update
            err ++ // # errors
        endif
    endfor
while ( err > 0 ) // repeat until no errors
return  $\mathbf{w}$ 
```

Pushes the correct  
sequence up and the  
incorrectly predicted one  
down

# Str. perceptron: algebraic interpretation

```
if ( $\mathbf{w} \cdot \varphi(\mathbf{x}^{(l)}, \hat{\mathbf{y}}) > \mathbf{w} \cdot \varphi(\mathbf{x}^{(l)}, \mathbf{y}^{(l)})$ )      // if mistake  
     $\mathbf{w} += \eta (\varphi(\mathbf{x}^{(l)}, \mathbf{y}^{(l)}) - \varphi(\mathbf{x}^{(l)}, \hat{\mathbf{y}}))$     // update  
endif
```

- ▶ We want after the update to increase  $\mathbf{w} \cdot (\varphi(\mathbf{x}^{(l)}, \mathbf{y}^{(l)}) - \varphi(\mathbf{x}^{(l)}, \hat{\mathbf{y}}))$ 
  - ▶ If the increase is large enough then  $\mathbf{y}^{(l)}$  will be scored above  $\hat{\mathbf{y}}$
- ▶ Clearly, that this is achieved as this product will be increased by

$$\eta \|\varphi(\mathbf{x}^{(l)}, \mathbf{y}^{(l)}) - \varphi(\mathbf{x}^{(l)}, \hat{\mathbf{y}})\|^2$$

*There might be other  
 $\mathbf{y}' \in \mathcal{Y}(\mathbf{x}^{(l)})$   
but we will deal with them  
on the next iterations*

# Structured Perceptron

---

## ▶ Positive:

- ▶ Very easy to implement
- ▶ Often, achieves respectable results
- ▶ As other discriminative techniques, does not make assumptions about the generative process
- ▶ Additional features can be easily integrated, as long as decoding is tractable

## ▶ Drawbacks

- ▶ “Good” discriminative algorithms should optimize some measure which is closely related to the expected testing results: what perceptron is doing on non-linearly separable data seems not clear
- ▶ However, for the averaged (voted) version a generalization bound which generalization properties of Perceptron (Freund & Shapire 98)

## ▶ Later, we will consider more advance learning algorithms



# Outline

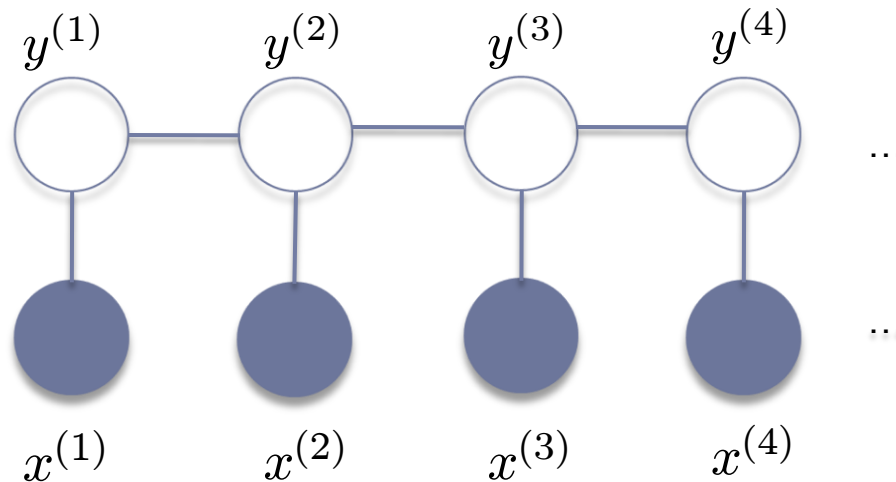
---

- ▶ Sequence labeling / segmentation problems: settings and example problems:
  - ▶ Part-of-speech tagging, named entity recognition, gesture recognition
- ▶ Hidden Markov Model
  - ▶ Standard definition + maximum likelihood estimation
  - ▶ General views: as a representative of linear models
- ▶ Perceptron and Structured Perceptron
  - ▶ algorithms / motivations
- ▶ Decoding with the Linear Model
- ▶ Discussion: Discriminative vs. Generative

# Decoding with the Linear model

---

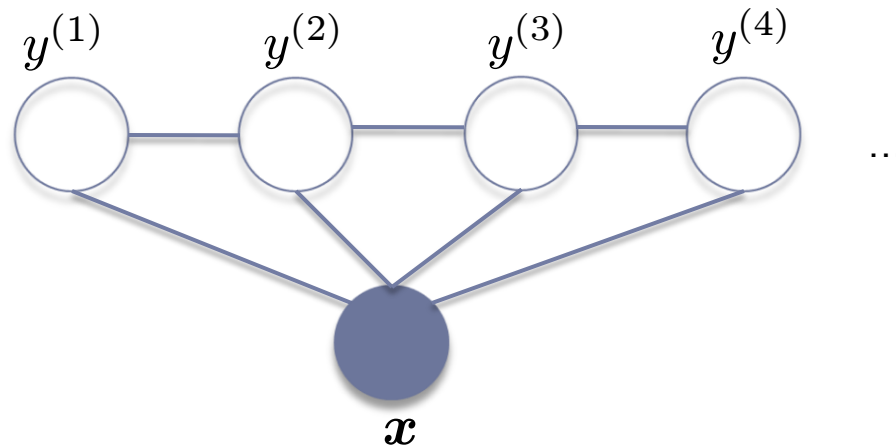
- ▶ Decoding:  $\mathbf{y} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} w \cdot \varphi(\mathbf{x}, \mathbf{y}')$
- ▶ Again a linear model with the following edge features (a generalization of a HMM)
- ▶ In fact, the algorithm does not depend on the feature of input (they do not need to be *local*)



# Decoding with the Linear model

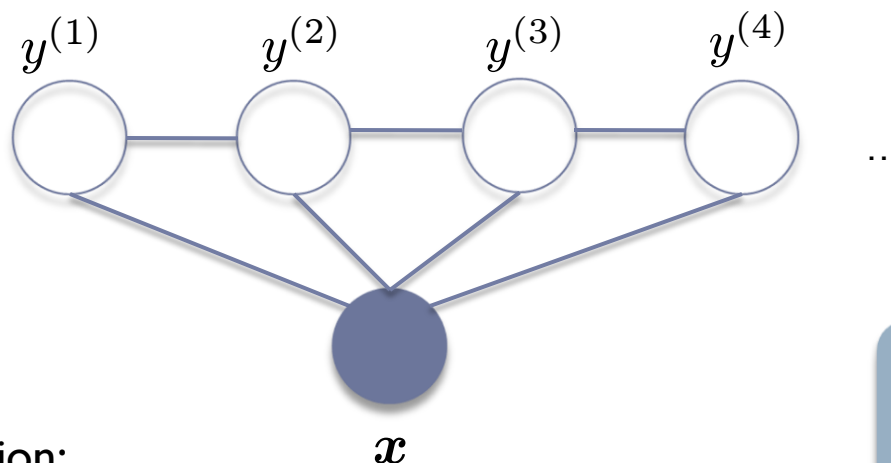
---

- ▶ Decoding:  $\mathbf{y} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} w \cdot \varphi(\mathbf{x}, \mathbf{y}')$
- ▶ Again a linear model with the following edge features (a generalization of a HMM)
- ▶ In fact, the algorithm does not depend on the feature of input (they do not need to be *local*)



# Decoding with the Linear model

- ▶ Decoding:  $\mathbf{y} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} w \cdot \varphi(\mathbf{x}, \mathbf{y}')$



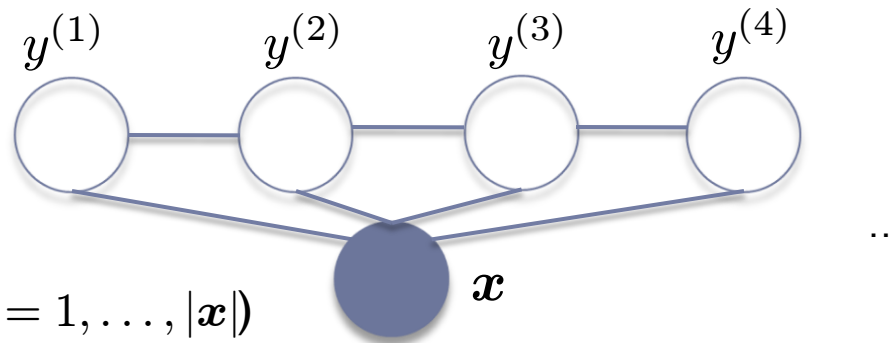
- ▶ Let's change notation:

Start/Stop symbol information (\$) can be encoded with them too.

- ▶ Edge scores  $f_t(y_{t-1}, y_t, \mathbf{x})$ : roughly corresponds to  $\log a_{y_{t-1}, y_t} + \log b_{y_t, \mathbf{x}_t}$
- ▶ Defined for  $t = 0$  too ("start" feature:  $y_0 = \$$  )
- ▶ Decode:  $\mathbf{y} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} \sum_{t=1}^{|\mathbf{x}|} f_t(y'_{t-1}, y'_t, \mathbf{x})$
- ▶ Decoding: a dynamic programming algorithm - Viterbi algorithm

# Viterbi algorithm

- ▶ **Decoding:**  $\mathbf{y} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{x})} \sum_{t=1}^{|\mathbf{x}|} f_t(y'_{t-1}, y'_t, \mathbf{x})$



- ▶ Loop invariant: ( $t = 1, \dots, |\mathbf{x}|$ )
  - ▶  $\text{score}_t[y]$  - score of the highest scoring sequence up to position  $t$  with
  - ▶  $\text{prev}_t[y]$  - previous tag on this sequence
- ▶ Init:  $\text{score}_0[\$] = 0$ ,  $\text{score}_0[y] = -\infty$  for other  $y$
- ▶ Recomputation ( $t = 1, \dots, |\mathbf{x}|$ )
$$\text{prev}_t[y] = \operatorname{argmax}_{y'} \text{score}_t[y'] + f_t(y', y, \mathbf{x})$$
$$\text{score}_t[y] = \text{score}_{t-1}[\text{prev}_t[y]] + f_t(\text{prev}_t[y], y, \mathbf{x})$$
- ▶ Return: retrace prev pointers starting from  $\operatorname{argmax}_y \text{score}_{|\mathbf{x}|}[y]$

Time complexity?  
 $O(N^2|\mathbf{x}|)$

# Outline

---

- ▶ Sequence labeling / segmentation problems: settings and example problems:
  - ▶ Part-of-speech tagging, named entity recognition, gesture recognition
- ▶ Hidden Markov Model
  - ▶ Standard definition + maximum likelihood estimation
  - ▶ General views: as a representative of linear models
- ▶ Perceptron and Structured Perceptron
  - ▶ algorithms / motivations
- ▶ Decoding with the Linear Model
- ▶ Discussion: Discriminative vs. Generative

# Recap: Sequence Labeling

---

- ▶ Hidden Markov Models:
  - ▶ How to estimate
- ▶ Discriminative models
  - ▶ How to learn with structured perceptron
- ▶ Both learning algorithms result in a linear model
  - ▶ How to label with the linear models

# Discriminative vs Generative

---

Not necessary the case  
for generative models  
with latent variables

## ▶ Generative models:

- ▶ Cheap to estimate: simply normalized counts
- ▶ Hard to integrate complex features: need to come up with a generative story and this story may be wrong
- ▶ Does not result in an optimal classifier when model assumptions are wrong (i.e., always)

## ▶ Discriminative models

- ▶ More expensive to learn: need to run decoding (here, Viterbi) during training and usually multiple times per an example
- ▶ Easy to integrate features: though some feature may make decoding intractable
- ▶ Usually less accurate on small datasets



# Reminders

---

- ▶ Speakers: slides about a week before the talk, meetings with me before/after this point will normally be needed
- ▶ Reviewers: reviews are accepted only before the day we consider the topic
- ▶ Everyone: References to the papers to read at GoogleDocs,
- ▶ These slides (and previous ones) will be online today
  - ▶ speakers: send me the last version of your slides too
- ▶ Next time: Lea about Models of Parsing, PCFGs vs general WCFGs (Michael Collins' book chapter)