

Predicting Structures: Conditional Models and Local Classifiers

CS 6355: Structured Prediction



Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences

Today's Agenda

- Conditional models for predicting sequences
- Maximum Entropy Markov Models

Today's Agenda

- Conditional models for predicting sequences
- Maximum Entropy Markov Models

HMM redux

- The independence assumption

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1} \mid y_i) \prod_{i=1}^n P(x_i \mid y_i)$$

HMM redux

- The independence assumption

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1} | y_i) \prod_{i=1}^n P(x_i | y_i)$$

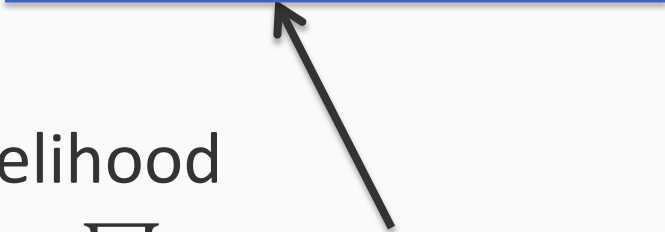
- Training via maximum likelihood

$$\max_{\pi, A, B} P(D | \pi, A, B) = \max_{\pi, A, B} \prod_i P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

We are optimizing joint likelihood of the input and the output for training

HMM redux

- The independence assumption

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1} | y_i) \prod_{i=1}^n P(x_i | y_i)$$


- Training via maximum likelihood

$$\max_{\pi, A, B} P(D | \pi, A, B) = \max_{\pi, A, B} \prod_i P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

We are optimizing joint likelihood of the input and the output for training

HMM redux

- The independence assumption

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1} | y_i) \prod_{i=1}^n P(x_i | y_i)$$

Probability of
input given the
prediction!

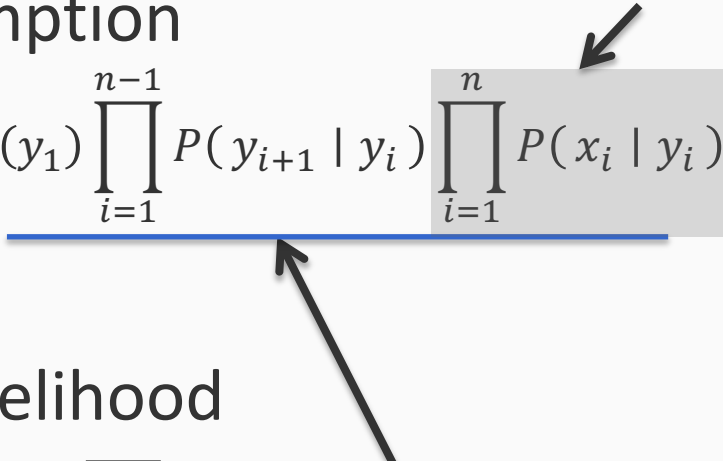
- Training via maximum likelihood

$$\max_{\pi, A, B} P(D | \pi, A, B) = \max_{\pi, A, B} \prod_i P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

We are optimizing joint likelihood of the input and the output for training

HMM redux

- The independence assumption

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1} | y_i) \prod_{i=1}^n P(x_i | y_i)$$


- Training via maximum likelihood

$$\max_{\pi, A, B} P(D | \pi, A, B) = \max_{\pi, A, B} \prod_i P(\mathbf{x}_i, y_i | \pi, A, B)$$

We are optimizing joint likelihood of the input and the output for training

At prediction time, we only care about the probability of output given the input:

$$P(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n)$$

Why not directly optimize this *conditional likelihood* instead?

Modeling next-state directly

- Instead of modeling the joint distribution $P(\mathbf{x}, \mathbf{y})$, focus on $P(\mathbf{y} \mid \mathbf{x})$ only
 - Which is what we care about eventually anyway
(At least in this context)

- For sequences, different formulations
 - Maximum Entropy Markov Model [McCallum, et al 2000]

(other names: discriminative/conditional Markov model, projection-based Markov model...)

Generative vs Discriminative models

- Generative models
 - learn $P(x, y)$
 - Characterize how the data is generated (both inputs and outputs)
 - Eg: Naïve Bayes, Hidden Markov Model
- Discriminative models
 - learn $P(y | x)$
 - Directly characterizes the decision boundary only
 - Eg: Logistic Regression, Conditional models (several names)

Generative vs Discriminative models

- Generative models
 - learn $P(x, y)$
 - Characterize how the data is generated (both inputs and outputs)
 - Eg: Naïve Bayes, Hidden Markov Model

A generative model tries to characterize the distribution of the inputs, a discriminative model doesn't care

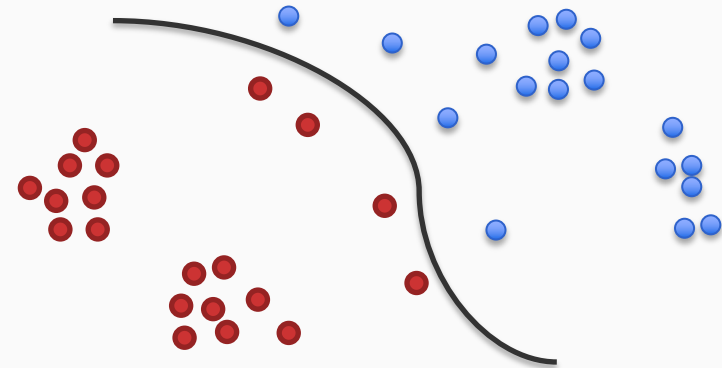
- Discriminative models
 - learn $P(y | x)$
 - Directly characterizes the decision boundary only
 - Eg: Logistic Regression, Conditional models (several names)

Generative vs Discriminative models

- Generative models
 - learn $P(x, y)$
 - Characterize how the data is generated (both inputs and outputs)
 - Eg: Naïve Bayes, Hidden Markov Model

A generative model tries to characterize the distribution of the inputs, a discriminative model doesn't care

- Discriminative models
 - learn $P(y | x)$
 - Directly characterizes the decision boundary only
 - Eg: Logistic Regression, Conditional models (several names)



Generative vs Discriminative models

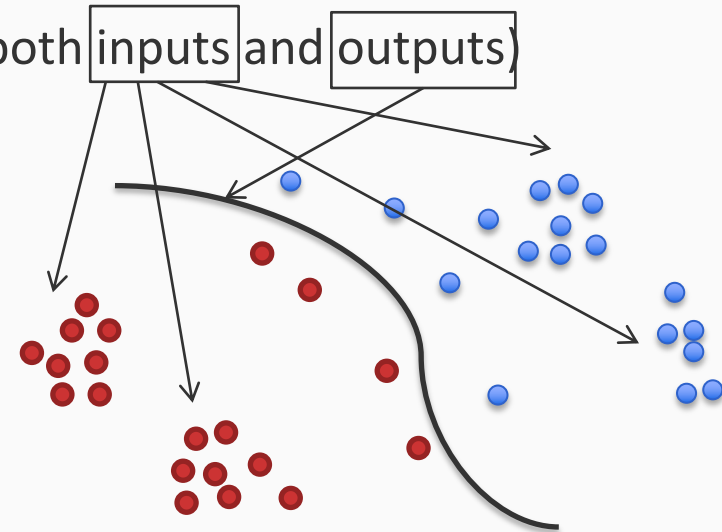
- Generative models

- learn $P(x, y)$
- Characterize how the data is generated (both **inputs** and **outputs**)
- Eg: Naïve Bayes, Hidden Markov Model

A generative model tries to characterize the distribution of the inputs, a discriminative model doesn't care

- Discriminative models

- learn $P(y | x)$
- Directly characterizes the decision boundary only
- Eg: Logistic Regression, Conditional models (several names)



Generative vs Discriminative models

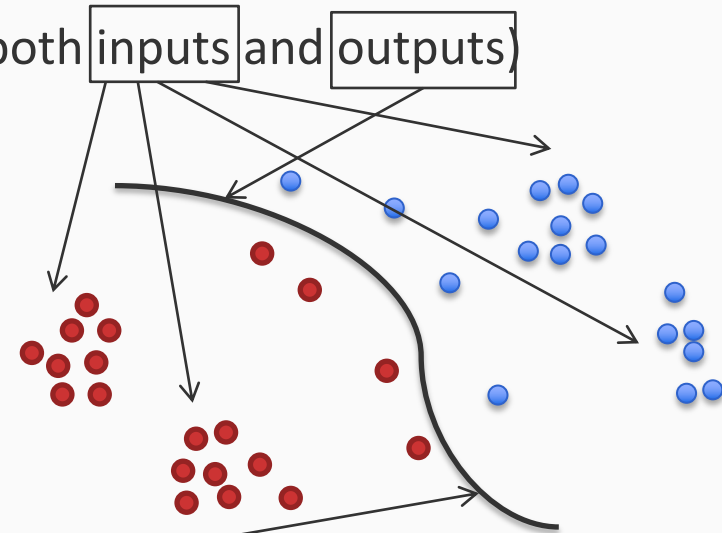
- Generative models

- learn $P(x, y)$
- Characterize how the data is generated (both **inputs** and **outputs**)
- Eg: Naïve Bayes, Hidden Markov Model

A generative model tries to characterize the distribution of the inputs, a discriminative model doesn't care

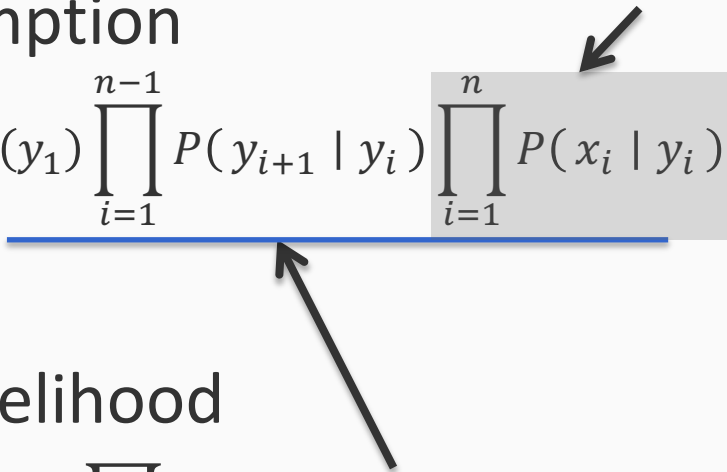
- Discriminative models

- learn $P(y | x)$
- Directly characterizes the **decision boundary** only
- Eg: Logistic Regression, Conditional models (several names)



HMM redux

- The independence assumption

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1} | y_i) \prod_{i=1}^n P(x_i | y_i)$$


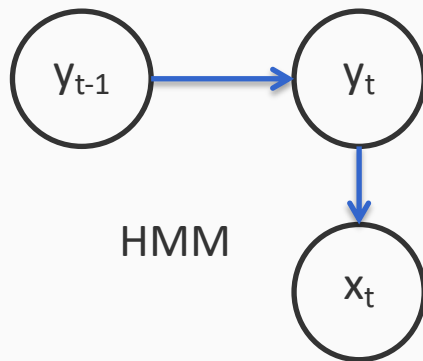
- Training via maximum likelihood

$$\max_{\pi, A, B} P(D | \pi, A, B) = \max_{\pi, A, B} \prod_i P(\mathbf{x}_i, y_i | \pi, A, B)$$

We are optimizing joint likelihood of the input and the output for training

At prediction time, we only care about the probability of output given the input. Why not directly optimize this *conditional likelihood* instead?

Let's revisit the independence assumptions

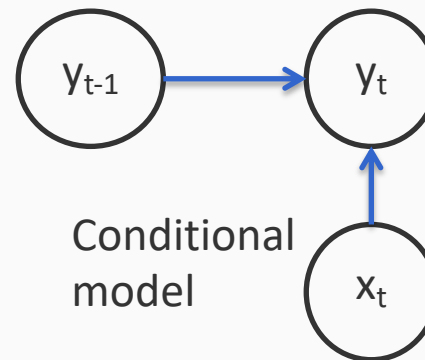
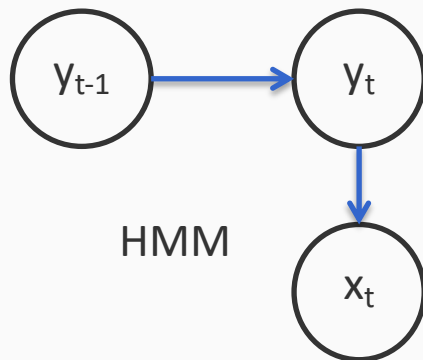


$$P(y_i | y_{i-1}, \text{anything else}) = P(y_i | y_{i-1})$$

$$P(x_i | y_i, \text{anything else}) = P(x_i | y_i)$$

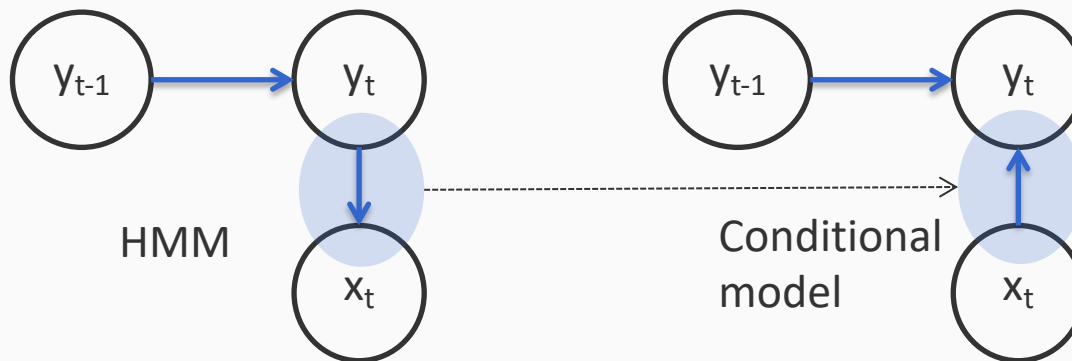
Another independence assumption

$$P(y_i \mid y_{i-1}, y_{i-2}, \dots, x_i, x_{i-1}, \dots) = P(y_i \mid y_{i-1}, x_i)$$



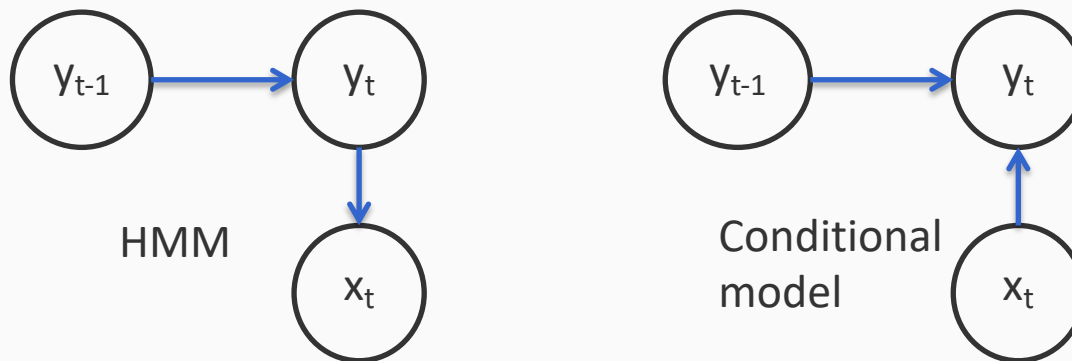
Another independence assumption

$$P(y_i \mid y_{i-1}, y_{i-2}, \dots, x_i, x_{i-1}, \dots) = P(y_i \mid y_{i-1}, x_i)$$



Another independence assumption

$$P(y_i \mid \textcolor{blue}{y}_{i-1}, y_{i-2}, \dots, \textcolor{blue}{x}_i, x_{i-1}, \dots) = P(y_i \mid \textcolor{blue}{y}_{i-1}, \textcolor{blue}{x}_i)$$

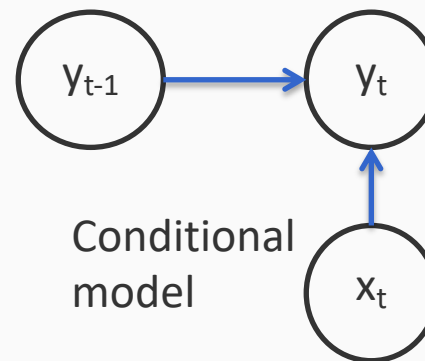
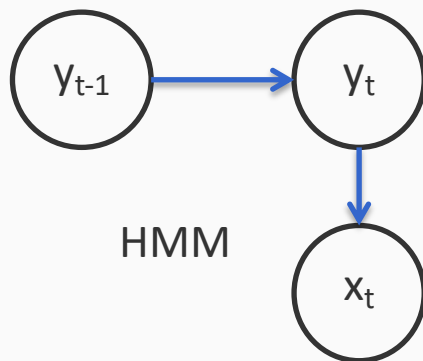


This assumption lets us write the conditional probability of the entire output sequence \mathbf{y} as

$$P(\mathbf{y} \mid \mathbf{x}) = \prod_i P(y_i \mid y_{i-1}, x_i)$$

Another independence assumption

$$P(y_i \mid y_{i-1}, y_{i-2}, \dots, x_i, x_{i-1}, \dots) = P(y_i \mid y_{i-1}, x_i)$$



This assumption lets us write the conditional probability of the entire output sequence \mathbf{y} as

$$P(\mathbf{y} \mid \mathbf{x}) = \prod_i P(y_i \mid y_{i-1}, x_i)$$

We need to learn this function

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

This is a multiclass classifier whose input is the pair y_{i-1}, x_i , and the label is the state y_i

Different approaches possible

1. Train a *maximum entropy* classifier (i.e., a multiclass logistic regression classifier)
2. Or, ignore the fact that we are predicting a probability, we only care about maximizing some *score*. Train any multiclass classifier, using say the perceptron algorithm

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

This is a multiclass classifier whose input is the pair y_{i-1}, x_i , and the label is the state y_i

Different approaches possible

1. Train a *maximum entropy* classifier (i.e., a multiclass logistic regression classifier)
2. Or, ignore the fact that we are predicting a probability, we only care about maximizing some *score*. Train any multiclass classifier, using say the perceptron algorithm

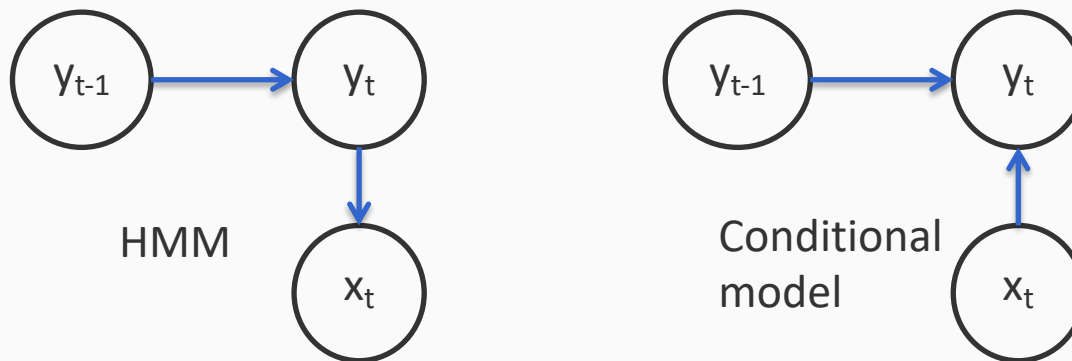
In any case, we can use any features from y_{i-1} and x_i . This was not possible with HMMs.

Today's Agenda

- Conditional models for predicting sequences
- Maximum Entropy Markov Models

The next-state model

$$P(y_i \mid y_{i-1}, y_{i-2}, \dots, x_i, x_{i-1}, \dots) = P(y_i \mid y_{i-1}, x_i)$$



This assumption lets us write the conditional probability of the output as

$$P(\mathbf{y}|\mathbf{x}) = \prod_i P(y_i|y_{i-1}, x_i)$$

We need to learn this function

Modeling ~~$P(y_i | y_{i-1}, x_i)$~~ $P(y_i | y_{i-1}, \mathbf{x})$

Different approaches possible

1. Train a *maximum entropy* classifier (i.e., a multiclass logistic regression classifier)
2. Or, ignore the fact that we are predicting a probability, we only care about maximizing some *score*. Train any multiclass classifier, using say the perceptron algorithm

Modeling ~~$P(y_i | y_{i-1}, x_i)$~~ $P(y_i | y_{i-1}, \mathbf{x})$

Different approaches possible

1. Train a *maximum entropy* classifier (i.e., a multiclass logistic regression classifier)
2. Or, ignore the fact that we are predicting a probability, we only care about maximizing some *score*. Train any multiclass classifier, using say the perceptron algorithm

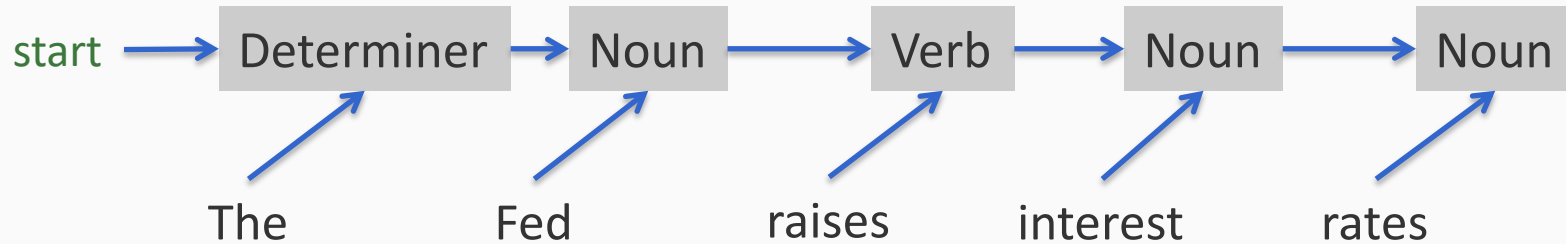
For both cases:

- Use rich features that depend on input and previous state
- We can increase the dependency to arbitrary neighboring x_i 's
 - Eg. Neighboring words influence this words POS tag

Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$

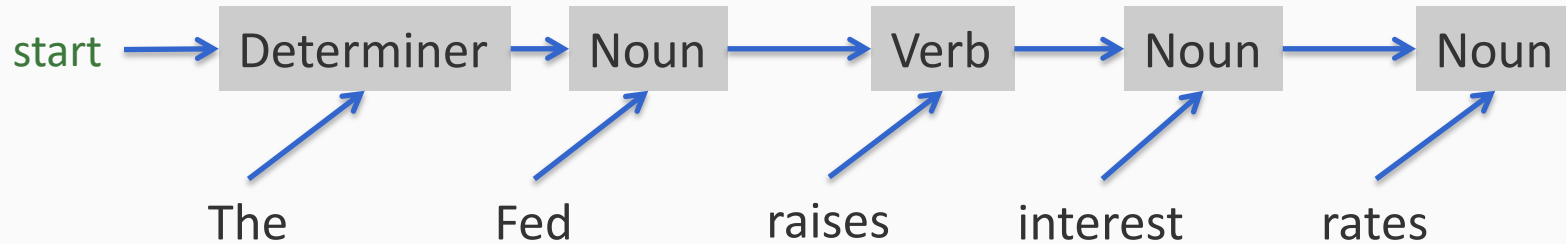


The **prediction task**: Using the entire input and the current label, predict the next label

Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$



word

Caps

-es

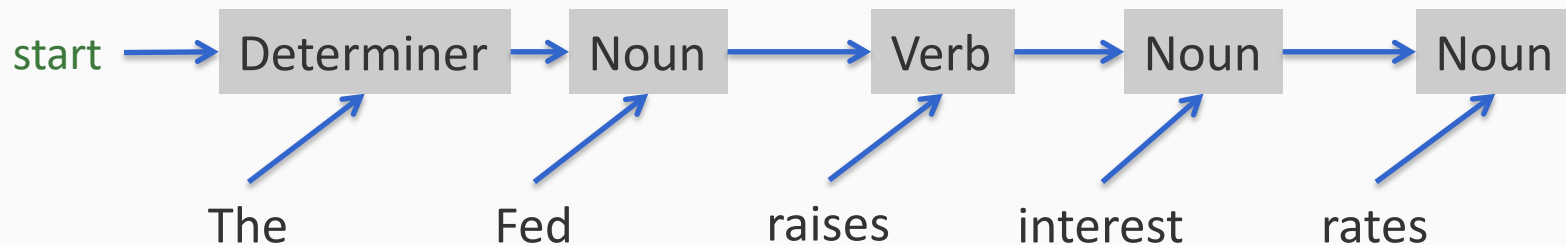
Previous

To model the probability, first, we need to define features for the current classification problem

Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$



word
Caps
-es?
Previous

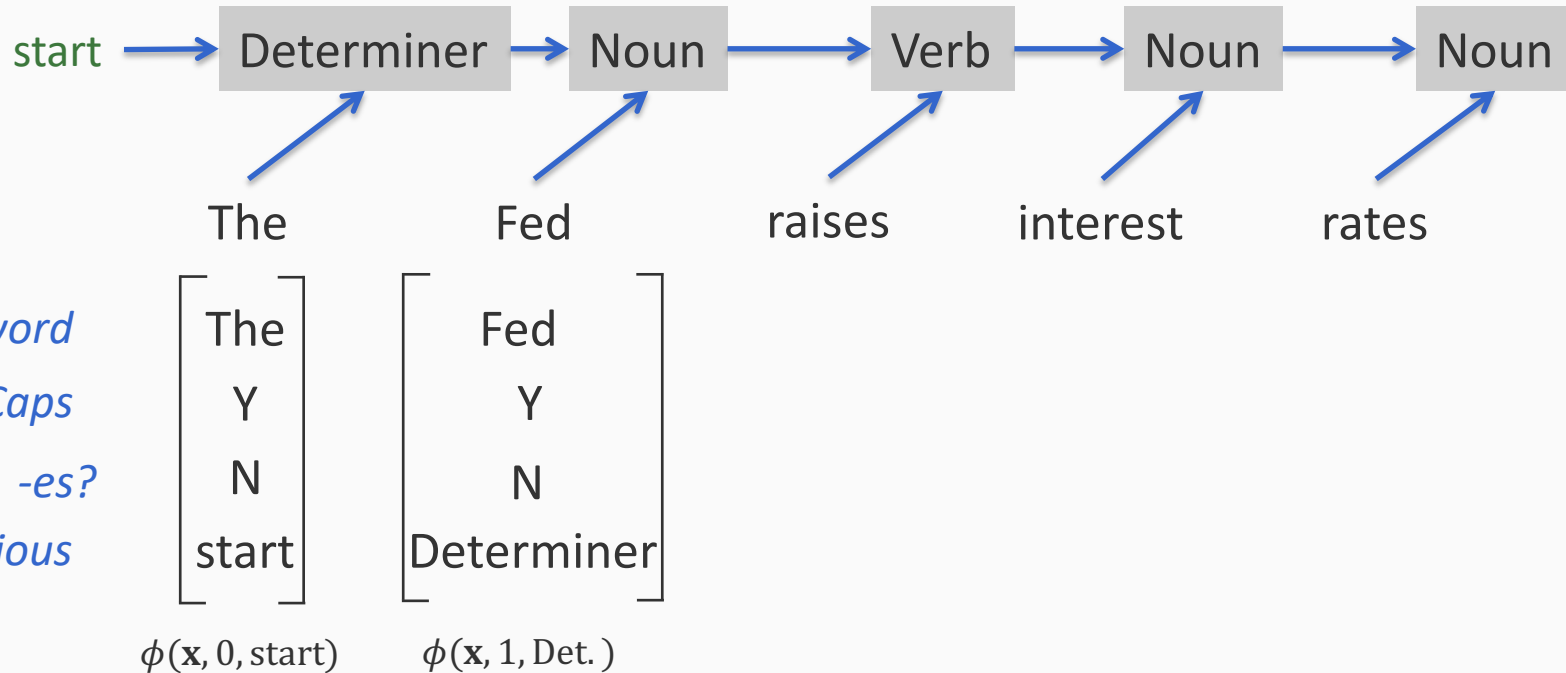
The
Y
N
start

$\phi(\mathbf{x}, 0, \text{start})$

Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

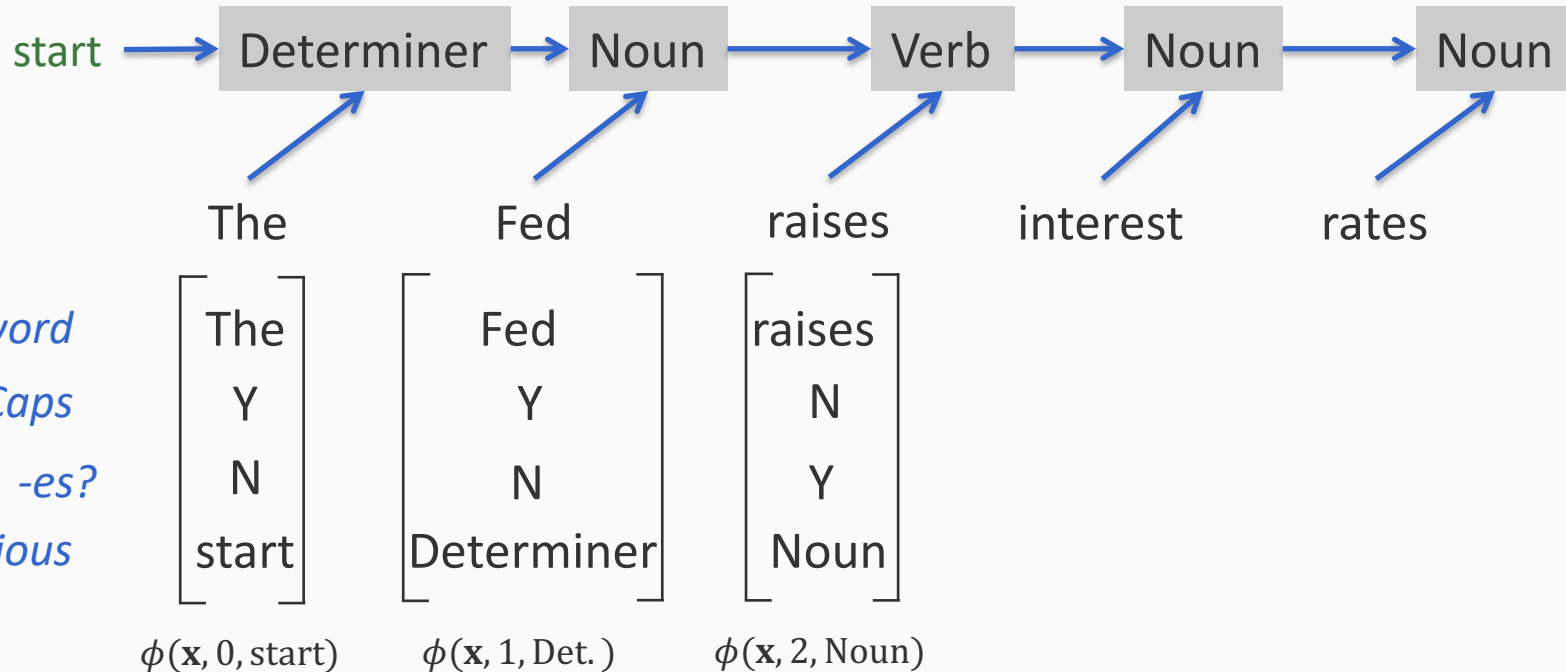
$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$



Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

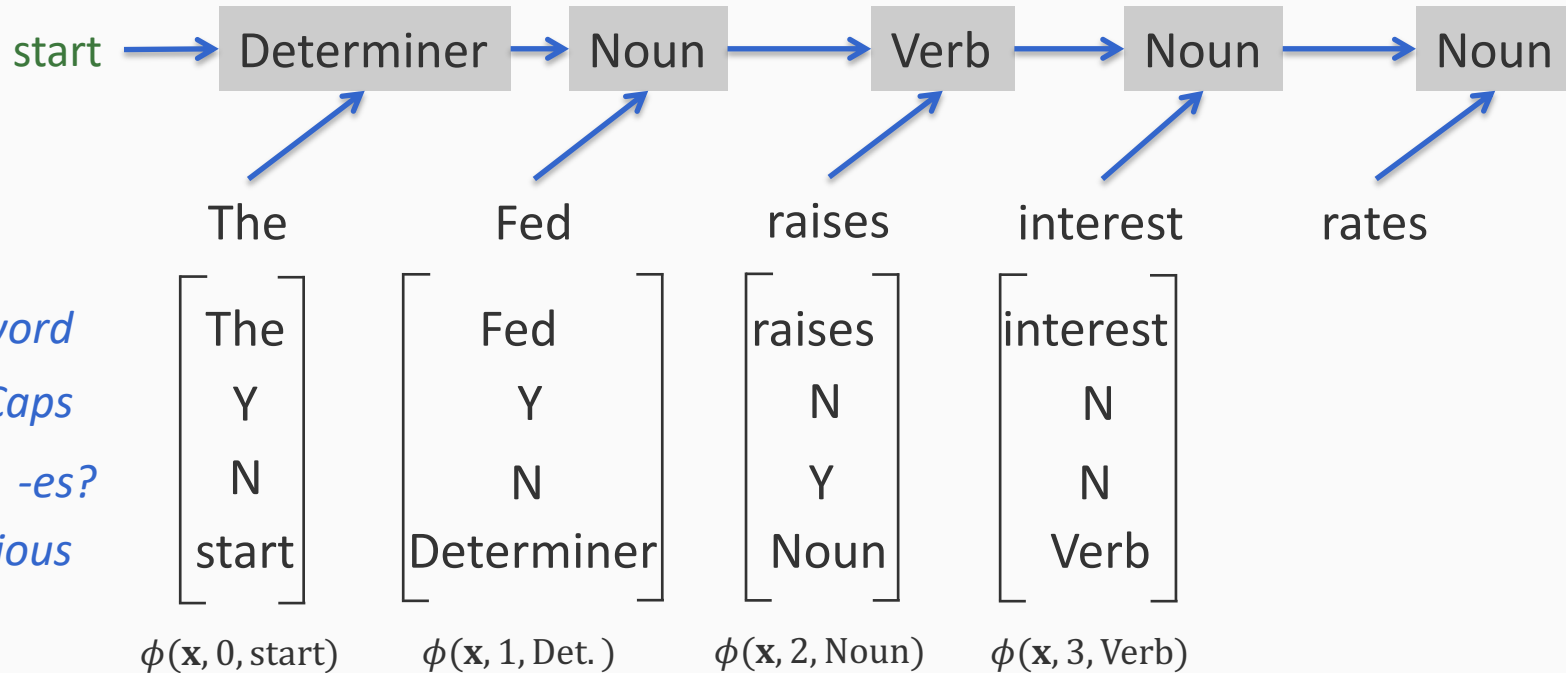
$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$



Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

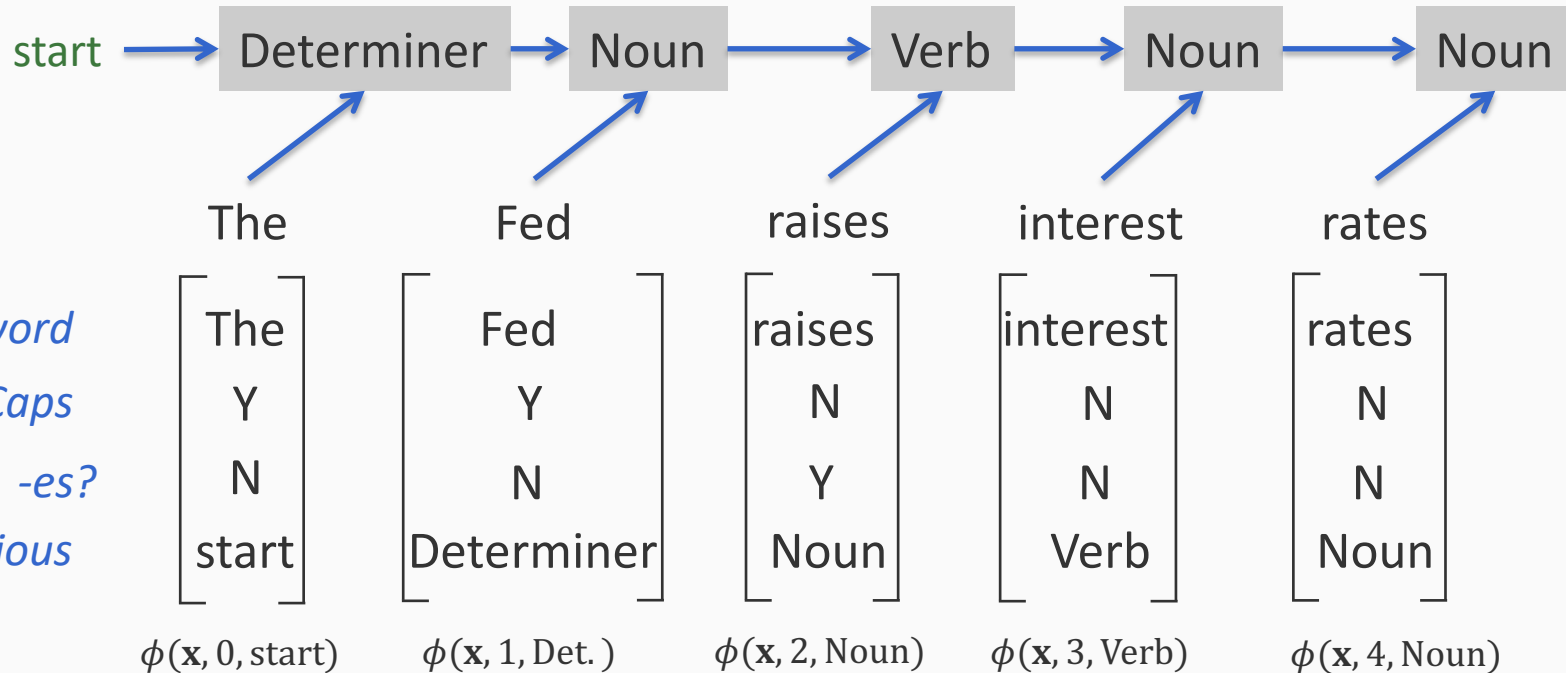
$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$



Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

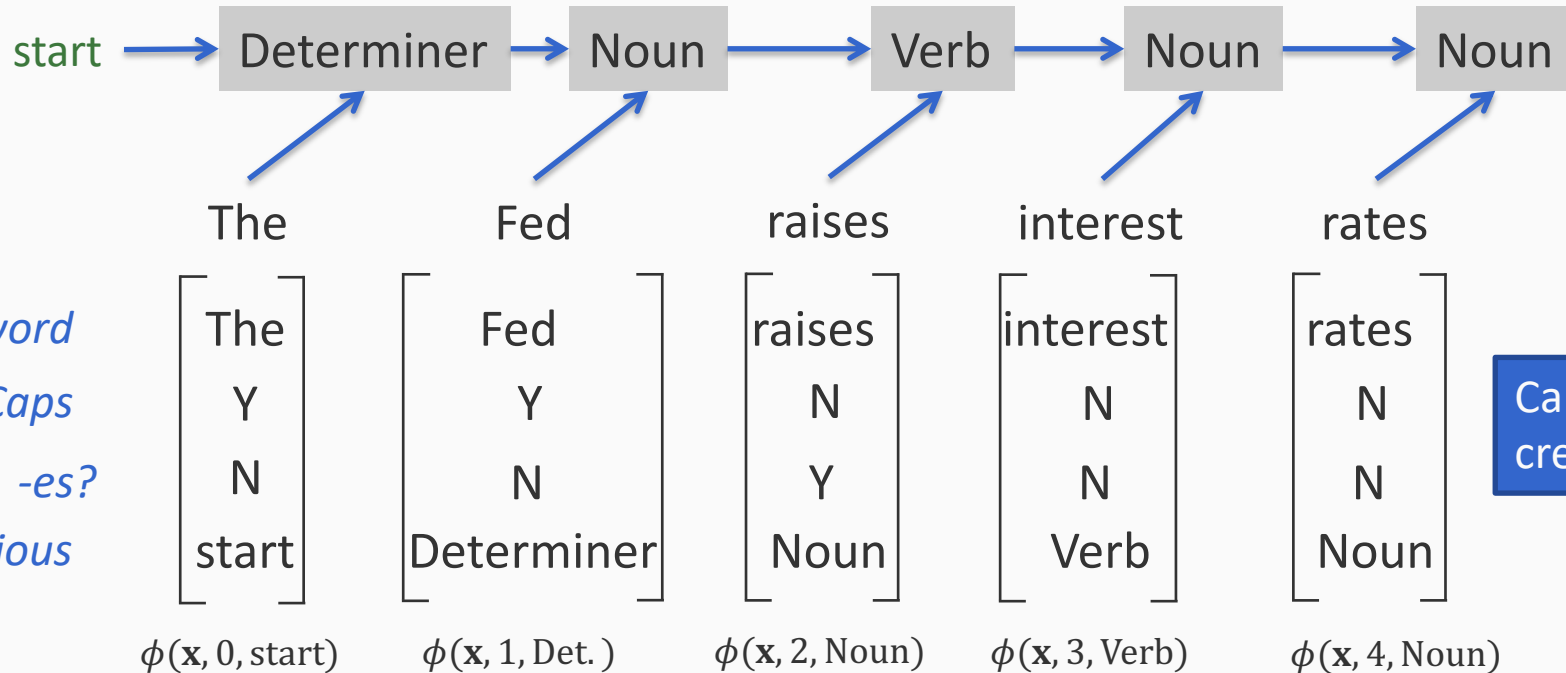
$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$



Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

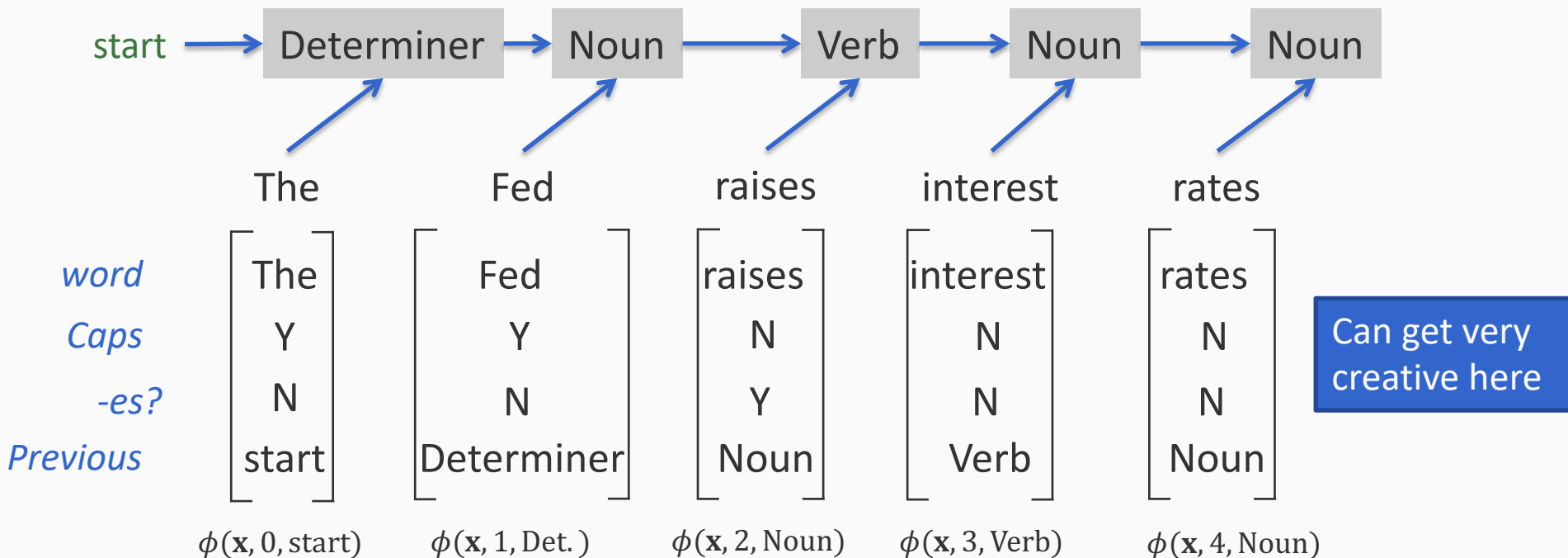
$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$



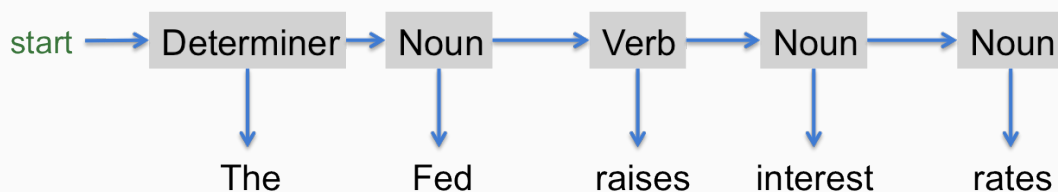
Maximum Entropy Markov Model

Goal: Compute $P(\mathbf{y} \mid \mathbf{x})$

$$P(y_i | y_{i-1}, \mathbf{x}) \propto \exp(\mathbf{w}^T \phi(\mathbf{x}, i, y_i, y_{i-1}))$$

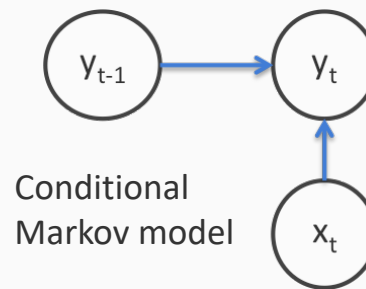
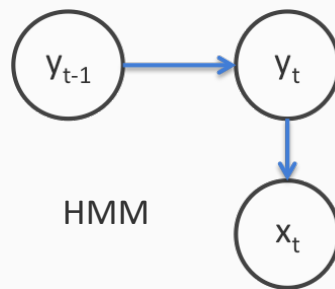


Compare to HMM: Only depends on the word and the previous tag



Using MEMM

- Training
 - Next-state predictor **locally** as maximum likelihood
 - Similar to any maximum entropy classifier
- Prediction/decoding
 - Modify the Viterbi algorithm for the new independence assumptions



$$\text{score}_i(s) = \max_{y_{i-1}} P(s|y_{i-1})P(x_i|s)\text{score}_{i-1}(y_{i-1})$$

$$\text{score}_i(s) = \max_{y_{i-1}} P(s|y_{i-1}, x_i)\text{score}_{i-1}(y_{i-1})$$

Generalization: Any multiclass classifier

- **Viterbi decoding**: we only need a score for each decision
 - So far, probabilistic classifiers
- In general, use any learning algorithm to build get a score for the label y_i given y_{i-1} and \mathbf{x}
 - Multiclass versions of perceptron, SVM
 - Just like MEMM, these allow arbitrary features to be defined

Exercise: Viterbi needs to be re-defined to work with sum of scores rather than the product of probabilities

Comparison to HMM

What we gain

1. Rich feature representation for inputs

- Helps generalize better by thinking about properties of the input tokens rather than the entire tokens
- Eg: If a word ends with –es, it might be a present tense verb (such as raises). Could be a feature; HMM cannot capture this

2. Discriminative predictor

- Model $P(\mathbf{y} \mid \mathbf{x})$ rather than $P(\mathbf{y}, \mathbf{x})$
- *Joint vs conditional*

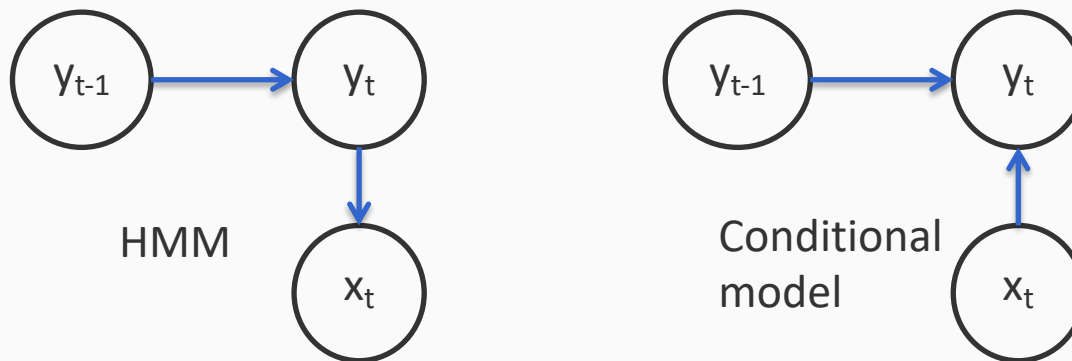
Questions?

Outline

- Conditional models for predicting sequences
- Log-linear models for multiclass classification
- Maximum Entropy Markov Models
 - The Label Bias Problem

The next-state model for sequences

$$P(y_i | \textcolor{red}{y}_{i-1}, y_{i-2}, \dots, \textcolor{red}{x}_i, x_{i-1}, \dots) = P(y_i | y_{i-1}, x_i)$$



This assumption lets us write the conditional probability of the output as

$$P(\mathbf{y}|\mathbf{x}) = \prod_i P(y_i | y_{i-1}, x_i)$$

We need to train local multiclass classifiers that predicts the next state given the previous state and the input

...local classifiers → Label bias problem

Let's look at the independence assumption

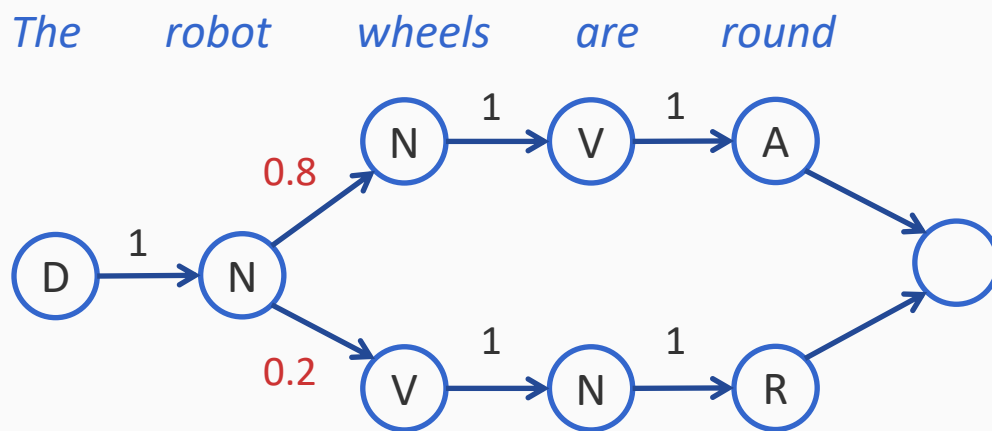
$$P(y_i | y_{i-1}, y_{i-2}, \dots, x_i, x_{i-1}, \dots) = P(y_i | y_{i-1}, x_i) \quad \text{"Next-state" classifiers are locally normalized}$$

...local classifiers → Label bias problem

Let's look at the independence assumption

$$P(y_i | \mathbf{y}_{i-1}, y_{i-2}, \dots, \mathbf{x}_i, x_{i-1}, \dots) = P(y_i | y_{i-1}, x_i) \quad \text{"Next-state" classifiers are locally normalized}$$

Eg: Part-of-speech tagging the sentence



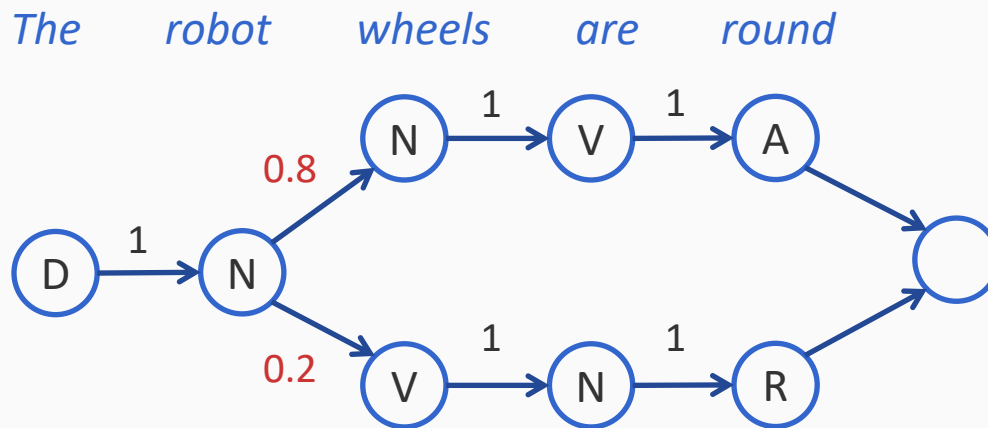
Suppose these are the only state transitions allowed

...local classifiers → Label bias problem

Let's look at the independence assumption

$$P(y_i | \mathbf{y}_{i-1}, y_{i-2}, \dots, \mathbf{x}_i, x_{i-1}, \dots) = P(y_i | y_{i-1}, x_i) \quad \text{"Next-state" classifiers are locally normalized}$$

Eg: Part-of-speech tagging the sentence



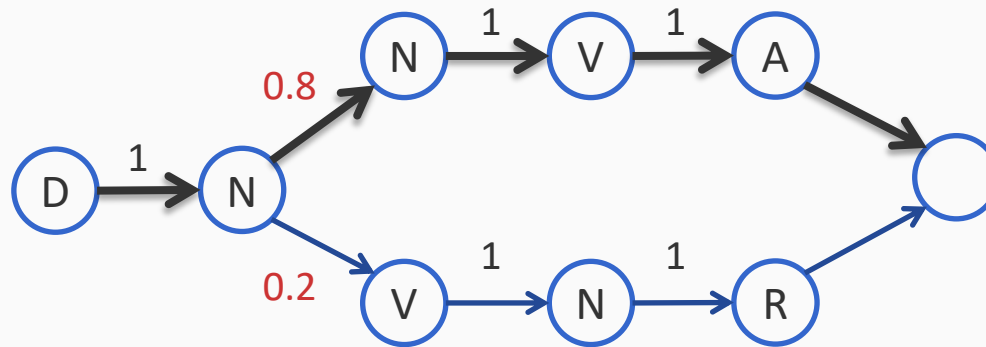
Option 1: $P(D \mid \text{The}) \times$
 $P(N \mid D, \text{robot}) \times$
 $P(N \mid N, \text{wheels}) \times$
 $P(V \mid N, \text{are}) \times$
 $P(A \mid V, \text{round})$

Option 2: $P(D \mid \text{The}) \times$
 $P(N \mid D, \text{robot}) \times$
 $P(V \mid N, \text{wheels}) \times$
 $P(N \mid V, \text{are}) \times$
 $P(R \mid N, \text{round})$

Suppose these are the only state transitions allowed

...local classifiers → Label bias problem

The robot wheels are round



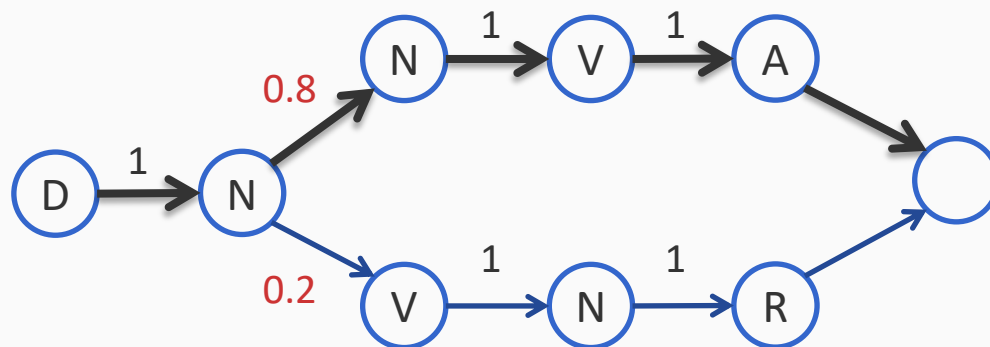
Suppose these are the only state transitions allowed

Option 1: $P(D \mid \text{The}) \times$
 $P(N \mid D, \text{robot}) \times$
 $P(N \mid N, \text{wheels}) \times$
 $P(V \mid N, \text{are}) \times$
 $P(A \mid V, \text{round})$

Option 2: $P(D \mid \text{The}) \times$
 $P(N \mid D, \text{robot}) \times$
 $P(V \mid N, \text{wheels}) \times$
 $P(N \mid V, \text{are}) \times$
 $P(R \mid N, \text{round})$

...local classifiers → Label bias problem

The robot wheels are round



Option 1: $P(D \mid \text{The}) \times$
 $P(N \mid D, \text{robot}) \times$
 $P(N \mid N, \text{wheels}) \times$
 ~~$P(V \mid N, \text{are})$~~ $\times P(V \mid N, \text{Fred}) \times$
 $P(A \mid V, \text{round})$

Option 2: $P(D \mid \text{The}) \times$
 $P(N \mid D, \text{robot}) \times$
 $P(V \mid N, \text{wheels}) \times$
 ~~$P(N \mid V, \text{are})$~~ $\times P(N \mid V, \text{Fred}) \times$
 $P(R \mid N, \text{round})$

Suppose these are the only state transitions allowed

The robot wheels Fred round

The path scores are the same

Even if the word Fred is never observed as a verb in the data, it will be predicted as one

The input *Fred* does not influence the output at all

Label Bias

- States with a single outgoing transition effectively ignore their input
 - States with lower-entropy next states are less influenced by observations
- Why?
 - Because each the next-state classifiers are locally normalized
 - If a state has fewer next states, each of those will get a higher probability mass
 - ...and hence preferred
- Side note: Surprisingly doesn't affect some tasks
 - Eg: part-of-speech tagging

Summary: Local models for Sequences

- Conditional models
- Use rich features in the mode
- Possibly suffer from label bias problem

(Other “local” models may have their own version of the label bias problem too.)

Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences

Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences