

Predicting Sequences: Structured Perceptron

CS 6355: Structured Prediction



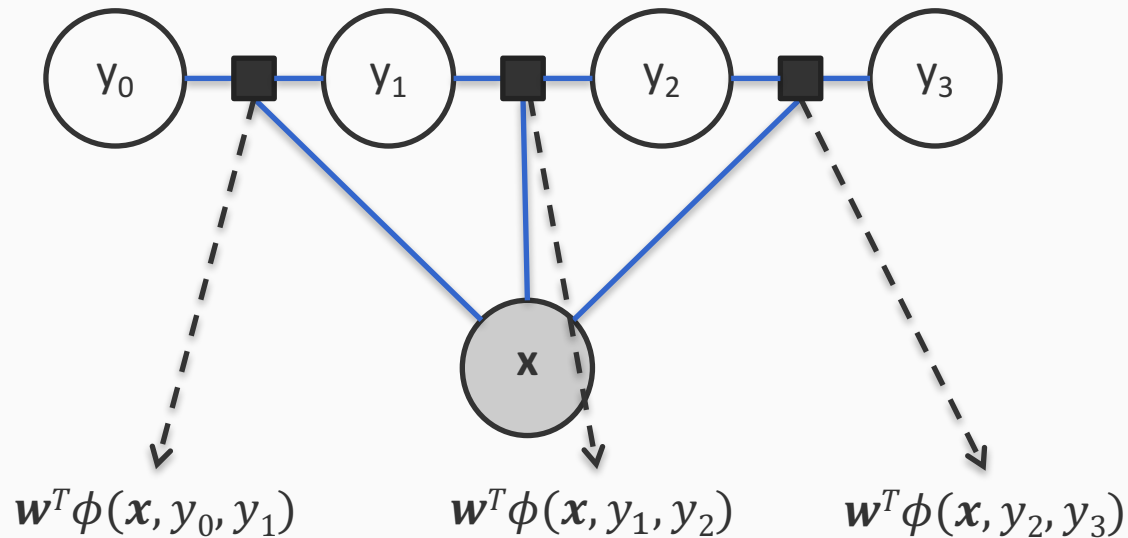
Conditional Random Fields summary

- An undirected graphical model
 - Decompose the score over the structure into a collection of factors
 - Each factor assigns a score to assignment of the random variables it is connected to
- Training and prediction
 - Final prediction via $\operatorname{argmax} w^T \phi(\mathbf{x}, \mathbf{y})$
 - Train by maximum (regularized) likelihood
- Connections to other models
 - Effectively a linear classifier
 - A generalization of logistic regression to structures
 - An conditional variant of a Markov Random Field
 - We will see this soon

Global features

The feature function decomposes over the sequence

$$\phi(\mathbf{x}, \mathbf{y}) = \sum_i \phi(\mathbf{x}, y_i, y_{i-1})$$



Outline

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

HMM is also a linear classifier

Consider the HMM:

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

HMM is also a linear classifier

Consider the HMM:

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

The diagram illustrates the decomposition of the HMM joint probability $P(\mathbf{x}, \mathbf{y})$ into a product of transition and emission probabilities. The equation is $P(\mathbf{x}, \mathbf{y}) = \prod_i P(y_i | y_{i-1}) P(x_i | y_i)$. The terms $P(y_i | y_{i-1})$ and $P(x_i | y_i)$ are highlighted in blue boxes. Dashed arrows point from these boxes to the labels "Transitions" and "Emissions" respectively.

HMM is also a linear classifier

Consider the HMM:

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

Or equivalently

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

Log joint probability = transition scores + emission scores

HMM is also a linear classifier

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

Log joint probability = transition scores + emission scores

Let us examine this expression using a carefully defined set of *indicator functions*

HMM is also a linear classifier

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

Log joint probability = transition scores + emission scores

Let us examine this expression using a carefully defined set of *indicator functions*

$$I_{[z]} = \begin{cases} 1, & z \text{ is true,} \\ 0, & z \text{ is false.} \end{cases}$$

Indicators are functions that map Booleans to 0 or 1

HMM is also a linear classifier

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

Log joint probability = transition scores + emission scores

Let us examine this expression using a carefully defined set of *indicator functions*

Equivalent to

$$\sum_s \sum_{s'} \log P(s | s') \cdot I_{[y_i=s]} \cdot I_{[y_{i-1}=s']}$$

The indicators ensure that only one of the elements of the double summation is non-zero

HMM is also a linear classifier

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

Log joint probability = transition scores + emission scores

Let us examine this expression using a carefully defined set of *indicator functions*

Equivalent to

$$\sum_s \log P(x_i | s) \cdot I_{[y_i=s]}$$

The indicators ensure that only one of the elements of the summation is non-zero

HMM is also a linear classifier

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

Let us examine this expression using a carefully defined set of *indicator functions*

$$\begin{aligned} \log P(\mathbf{x}, \mathbf{y}) = & \sum_i \sum_s \sum_{s'} \log P(s | s') \cdot I_{[y_i=s]} \cdot I_{[y_{i-1}=s']} \\ & + \sum_i \sum_s \log P(x_i | s) \cdot I_{[y_i=s]} \end{aligned}$$

HMM is also a linear classifier

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

Let us examine this expression using a carefully defined set of *indicator functions*

$$\begin{aligned} \log P(\mathbf{x}, \mathbf{y}) &= \sum_i \sum_s \sum_{s'} \log P(s | s') \cdot I_{[y_i=s]} \cdot I_{[y_{i-1}=s']} \\ &\quad + \sum_i \sum_s \log P(x_i | s) \cdot I_{[y_i=s]} \end{aligned}$$

$$\begin{aligned} \log P(\mathbf{x}, \mathbf{y}) &= \sum_s \sum_{s'} \log P(s | s') \sum_i I_{[y_i=s]} I_{[y_{i-1}=s']} \\ &\quad + \sum_s \log P(x_i | s) \sum_i I_{[y_i=s]} \end{aligned}$$

HMM is also a linear classifier

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

Let us examine this expression using a carefully defined set of *indicator functions*

$$\log P(\mathbf{x}, \mathbf{y}) = \sum_s \sum_{s'} \log P(s | s') \sum_i I_{[y_i=s]} I_{[y_{i-1}=s']} + \sum_s \log P(x_i | s) \sum_i I_{[y_i=s]}$$

Number of times
there is a transition in
the sequence from
state s' to state s
 $\text{Count}(s' \rightarrow s)$

HMM is also a linear classifier

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

Let us examine this expression using a carefully defined set of *indicator functions*

$$\begin{aligned} \log P(\mathbf{x}, \mathbf{y}) = & \sum_s \sum_{s'} \log P(s | s') \cdot \text{Count}(s' \rightarrow s) \\ & + \sum_s \log P(x_i | s) \sum_i I_{[y_i=s]} \end{aligned}$$

HMM is also a linear classifier

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

Let us examine this expression using a carefully defined set of *indicator functions*

$$\begin{aligned} \log P(\mathbf{x}, \mathbf{y}) = & \sum_s \sum_{s'} \log P(s | s') \cdot \text{Count}(s' \rightarrow s) \\ & + \sum_s \log P(x_i | s) \sum_i I_{[y_i=s]} \end{aligned}$$

Number of times state s occurs in the sequence: $\text{Count}(s)$

HMM is also a linear classifier

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

Let us examine this expression using a carefully defined set of *indicator functions*

$$\begin{aligned} \log P(\mathbf{x}, \mathbf{y}) = & \sum_s \sum_{s'} \log P(s | s') \cdot \text{Count}(s' \rightarrow s) \\ & + \sum_s \log P(x_i | s) \cdot \text{Count}(s) \end{aligned}$$

HMM is also a linear classifier

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

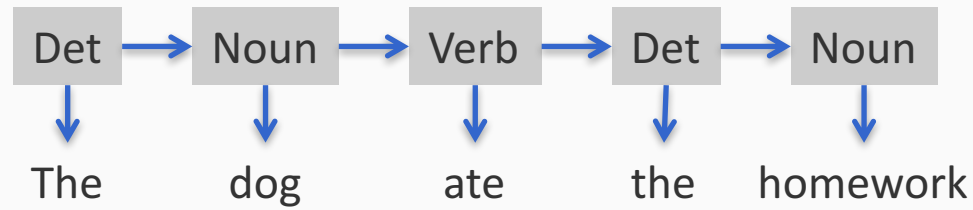
Let us examine this expression using a carefully defined set of *indicator functions*

$$\begin{aligned} \log P(\mathbf{x}, \mathbf{y}) &= \sum_s \sum_{s'} \log P(s | s') \cdot \text{Count}(s' \rightarrow s) \\ &\quad + \sum_s \log P(x_i | s) \cdot \text{Count}(s) \end{aligned}$$

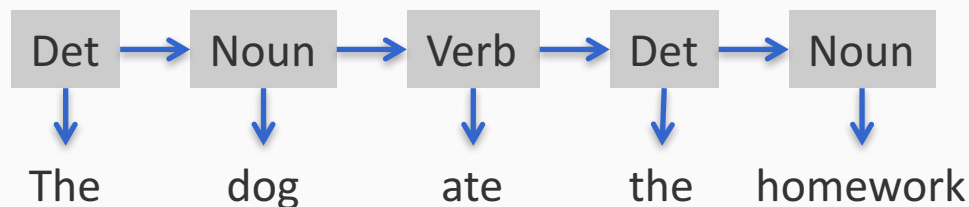
This is a linear function

log P terms are the weights; counts via indicators are features
Can be written as $\mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$ and add more features

HMM is a linear classifier: An example

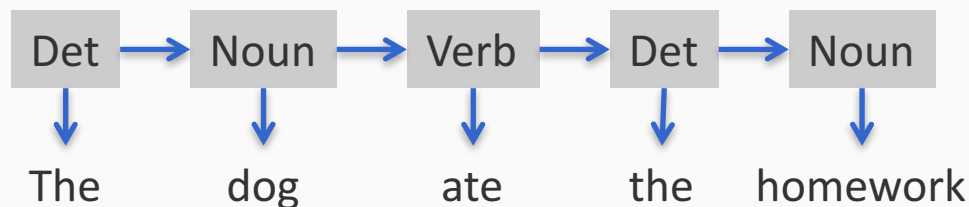


HMM is a linear classifier: An example



Consider $\log P(\mathbf{y}, \mathbf{x}) = \sum_i \log P(y_i | y_{i-1}) + \log P(x_i | y_i)$

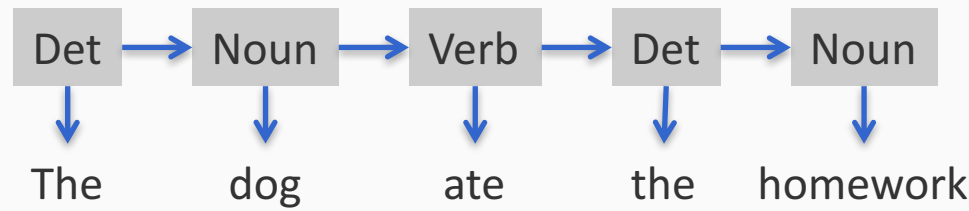
HMM is a linear classifier: An example



Consider $\log P(\mathbf{y}, \mathbf{x}) = \sum_i \log P(y_i | y_{i-1}) + \log P(x_i | y_i)$

Transition scores + Emission scores

HMM is a linear classifier: An example



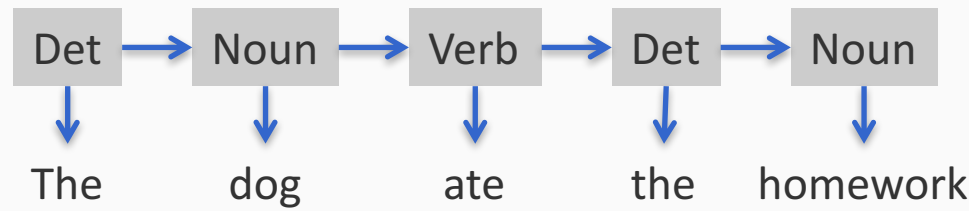
Consider $\log P(\mathbf{y}, \mathbf{x}) = \sum_i \log P(y_i | y_{i-1}) + \log P(x_i | y_i)$

$\log P(\text{Det} \rightarrow \text{Noun}) \times 2$

$+ \log P(\text{Noun} \rightarrow \text{Verb}) \times 1$ $+ \text{Emission scores}$

$+ \log P(\text{Verb} \rightarrow \text{Det}) \times 1$

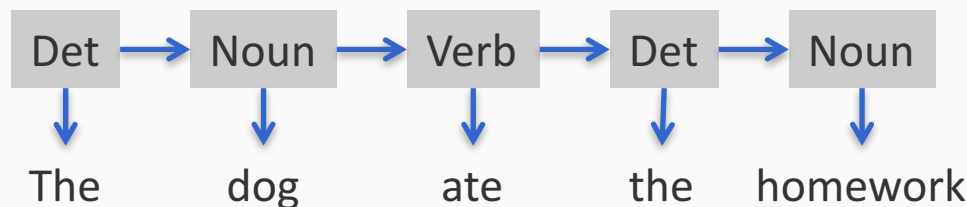
HMM is a linear classifier: An example



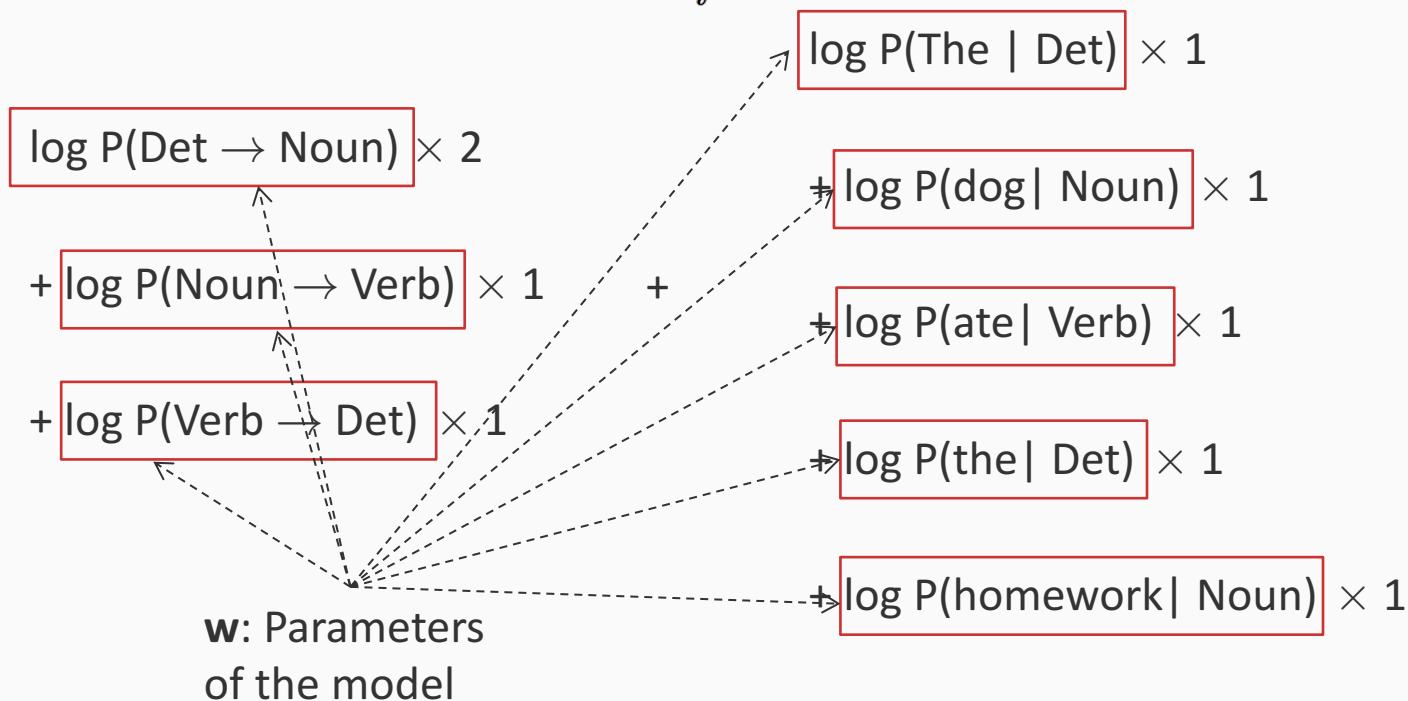
Consider $\log P(\mathbf{y}, \mathbf{x}) = \sum_i \log P(y_i | y_{i-1}) + \log P(x_i | y_i)$

$$\begin{aligned} & \log P(\text{The} | \text{Det}) \times 1 \\ \log P(\text{Det} \rightarrow \text{Noun}) \times 2 & + \log P(\text{dog} | \text{Noun}) \times 1 \\ + \log P(\text{Noun} \rightarrow \text{Verb}) \times 1 & + \\ + \log P(\text{Verb} \rightarrow \text{Det}) \times 1 & + \log P(\text{ate} | \text{Verb}) \times 1 \\ & + \log P(\text{the} | \text{Det}) \times 1 \\ & + \log P(\text{homework} | \text{Noun}) \times 1 \end{aligned}$$

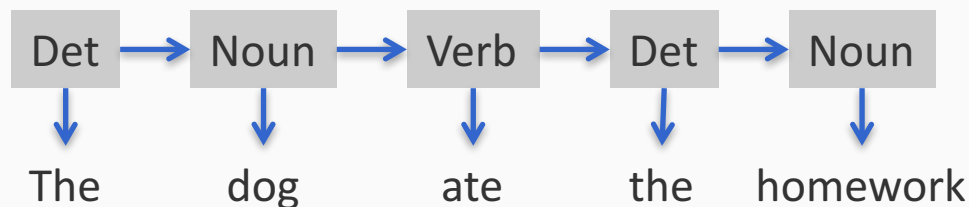
HMM is a linear classifier: An example



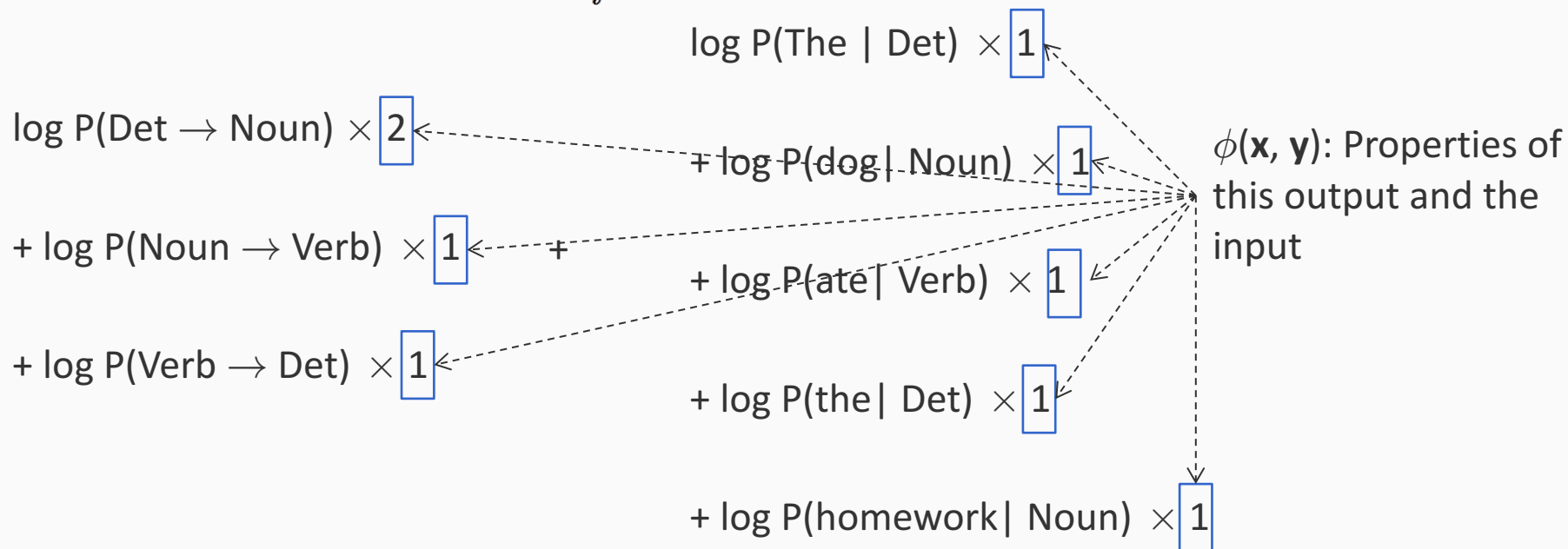
Consider $\log P(\mathbf{y}, \mathbf{x}) = \sum_i \log P(y_i | y_{i-1}) + \log P(x_i | y_i)$



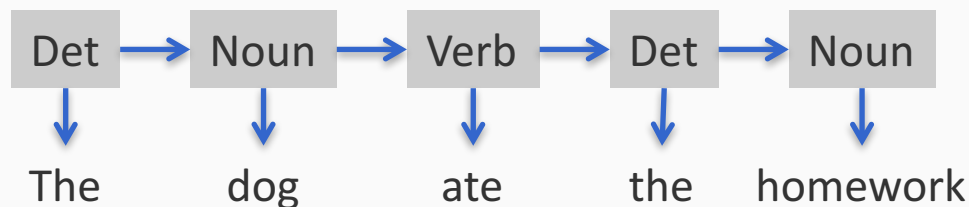
HMM is a linear classifier: An example



Consider $\log P(\mathbf{y}, \mathbf{x}) = \sum_i \log P(y_i | y_{i-1}) + \log P(x_i | y_i)$



HMM is a linear classifier: An example



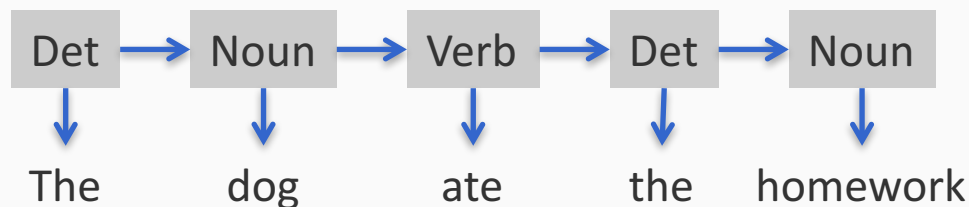
Consider $\log P(\mathbf{y}, \mathbf{x}) = \sum_i \log P(y_i | y_{i-1}) + \log P(x_i | y_i)$

$$\begin{bmatrix} \log P(Det \rightarrow Noun) \\ \log P(Noun \rightarrow Verb) \\ \log P(Verb \rightarrow Det) \\ \log P(The | Det) \\ \log P(dog | Noun) \\ \log P(ate | Verb) \\ \log P(the | Det) \\ \log P(homework | Noun) \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

\mathbf{w} : Parameters of the model

$\phi(\mathbf{x}, \mathbf{y})$: Properties of this output and the input

HMM is a linear classifier: An example



Consider $\log P(\mathbf{y}, \mathbf{x}) = \sum_i \log P(y_i | y_{i-1}) + \log P(x_i | y_i)$

$$\begin{bmatrix}
 \log P(Det \rightarrow Noun) \\
 \log P(Noun \rightarrow Verb) \\
 \log P(Verb \rightarrow Det) \\
 \log P(The | Det) \\
 \log P(dog | Noun) \\
 \log P(ate | Verb) \\
 \log P(the | Det) \\
 \log P(homework | Noun)
 \end{bmatrix}
 \cdot
 \begin{bmatrix}
 2 \\
 1 \\
 1 \\
 1 \\
 1 \\
 1 \\
 1 \\
 1 \\
 1
 \end{bmatrix}$$

$\log P(\mathbf{x}, \mathbf{y}) = \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$

\mathbf{w} : Parameters of the model

$\phi(\mathbf{x}, \mathbf{y})$: Properties of this output and the input

Towards structured Perceptron

1. HMM is a linear classifier

- Can we treat it as any linear classifier for training?
- If so, we could add additional features that are global properties
 - As long as the output can be decomposed for easy inference

Towards structured Perceptron

1. HMM is a linear classifier
 - Can we treat it as any linear classifier for training?
 - If so, we could add additional features that are global properties
 - As long as the output can be decomposed for easy inference
2. The Viterbi algorithm calculates $\max w^T \phi(\mathbf{x}, \mathbf{y})$

Viterbi only cares about scores to structures (not necessarily normalized)

Towards structured Perceptron

1. HMM is a linear classifier
 - Can we treat it as any linear classifier for training?
 - If so, we could add additional features that are global properties
 - As long as the output can be decomposed for easy inference
2. The Viterbi algorithm calculates $\max w^T \phi(\mathbf{x}, \mathbf{y})$

Viterbi only cares about scores to structures (not necessarily normalized)
3. We could push the learning algorithm to train for un-normalized scores
 - If we need normalization, we could always normalize by exponentiating and dividing by Z (the partition function)
 - That is, the learning algorithm can effectively just focus on the score of \mathbf{y} for a particular \mathbf{x}
 - *Train a discriminative model instead of the generative one!*

Structured Perceptron algorithm

Given a training set $D = \{(\mathbf{x}, \mathbf{y})\}$

1. Initialize $\mathbf{w} = \mathbf{0} \in \mathcal{R}^n$
2. For epoch = 1 ... T:
 1. For each training example $(\mathbf{x}, \mathbf{y}) \in D$:

Structured Perceptron update

3. Return \mathbf{w}

Prediction: $\operatorname{argmax}_{\mathbf{y}} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$

Structured Perceptron algorithm

Given a training set $D = \{(\mathbf{x}, \mathbf{y})\}$

1. Initialize $\mathbf{w} = \mathbf{0} \in \mathcal{R}^n$
2. For epoch = 1 ... T:
 1. For each training example $(\mathbf{x}, \mathbf{y}) \in D$:
 1. Predict $\mathbf{y}' = \operatorname{argmax}_{\mathbf{y}'} \mathbf{w}^\top \phi(\mathbf{x}, \mathbf{y}')$
 2. If $\mathbf{y} \neq \mathbf{y}'$, update $\mathbf{w} \leftarrow \mathbf{w} + r (\phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \mathbf{y}'))$
3. Return \mathbf{w}


Prediction: $\operatorname{argmax}_{\mathbf{y}} \mathbf{w}^\top \phi(\mathbf{x}, \mathbf{y})$

Structured Perceptron algorithm

Given a training set $D = \{(\mathbf{x}, \mathbf{y})\}$

1. Initialize $\mathbf{w} = \mathbf{0} \in \mathfrak{R}^n$
2. For epoch = 1 ... T:
 1. For each training example $(\mathbf{x}, \mathbf{y}) \in D$:
 1. Predict $\mathbf{y}' = \operatorname{argmax}_{\mathbf{y}'} \mathbf{w}^\top \phi(\mathbf{x}, \mathbf{y}')$
 2. If $\mathbf{y} \neq \mathbf{y}'$, update $\mathbf{w} \leftarrow \mathbf{w} + r (\phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \mathbf{y}'))$
3. Return \mathbf{w}

Prediction: $\operatorname{argmax}_{\mathbf{y}} \mathbf{w}^\top \phi(\mathbf{x}, \mathbf{y})$



Update only on an error.
Perceptron is a mistake-driven algorithm.
If there is a mistake, promote \mathbf{y} and demote \mathbf{y}'

Structured Perceptron algorithm

Given a training set $D = \{(\mathbf{x}, \mathbf{y})\}$

1. Initialize $\mathbf{w} = \mathbf{0} \in \mathcal{R}^n$

2. For epoch = 1 ... T:

T is a hyperparameter to the algorithm

1. For each training example $(\mathbf{x}, \mathbf{y}) \in D$:

1. Predict $\mathbf{y}' = \operatorname{argmax}_{\mathbf{y}'} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y}')$

2. If $\mathbf{y} \neq \mathbf{y}'$, update $\mathbf{w} \leftarrow \mathbf{w} + r (\phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \mathbf{y}'))$

3. Return \mathbf{w}

Prediction: $\operatorname{argmax}_{\mathbf{y}} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$

Structured Perceptron algorithm

Given a training set $D = \{(\mathbf{x}, \mathbf{y})\}$

1. Initialize $\mathbf{w} = \mathbf{0} \in \mathcal{R}^n$
2. For epoch = 1 ... T:
 1. For each training example $(\mathbf{x}, \mathbf{y}) \in D$:
 1. Predict $\mathbf{y}' = \operatorname{argmax}_{\mathbf{y}'} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y}')$
 2. If $\mathbf{y} \neq \mathbf{y}'$, update $\mathbf{w} \leftarrow \mathbf{w} + r (\phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \mathbf{y}'))$
3. Return \mathbf{w}

In practice, good to shuffle D before the inner loop

Prediction: $\operatorname{argmax}_{\mathbf{y}} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$

Structured Perceptron algorithm

Given a training set $D = \{(\mathbf{x}, \mathbf{y})\}$

1. Initialize $\mathbf{w} = \mathbf{0} \in \mathcal{R}^n$
2. For epoch = 1 ... T:
 1. For each training example $(\mathbf{x}, \mathbf{y}) \in D$:
 1. Predict $\mathbf{y}' = \operatorname{argmax}_{\mathbf{y}'} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y}')$
 2. If $\mathbf{y} \neq \mathbf{y}'$, update $\mathbf{w} \leftarrow \mathbf{w} + r (\phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \mathbf{y}'))$
3. Return \mathbf{w}

Inference in training loop!

Prediction: $\operatorname{argmax}_{\mathbf{y}} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$

Notes on structured perceptron

- Mistake bound for separable data, just like perceptron
- In practice, use averaging for better generalization
 - Initialize $\mathbf{a} = 0$
 - After each step, whether there is an update or not, $\mathbf{a} \leftarrow \mathbf{w} + \mathbf{a}$
 - Note, we still check for mistake using \mathbf{w} not \mathbf{a}
 - Return \mathbf{a} at the end instead of \mathbf{w}
 - **Exercise:** Optimize this for performance – modify \mathbf{a} *only* on errors
- Global update
 - One weight vector for entire sequence
 - Not for each position
 - Same algorithm can be derived via constraint classification
 - Create a binary classification data set and run perceptron

Structured Perceptron with averaging

Given a training set $D = \{(\mathbf{x}, \mathbf{y})\}$

1. Initialize $\mathbf{w} = \mathbf{0} \in \mathcal{R}^n$, $\mathbf{a} = \mathbf{0} \in \mathcal{R}^n$
2. For epoch = 1 ... T:
 1. For each training example $(\mathbf{x}, \mathbf{y}) \in D$:
 1. Predict $\mathbf{y}' = \operatorname{argmax}_{\mathbf{y}'} \mathbf{w}^\top \phi(\mathbf{x}, \mathbf{y}')$
 2. If $\mathbf{y} \neq \mathbf{y}'$, update $\mathbf{w} \leftarrow \mathbf{w} + r (\phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \mathbf{y}'))$
 3. Set $\mathbf{a} \leftarrow \mathbf{a} + \mathbf{w}$
3. Return \mathbf{a}

CRF vs. structured perceptron

Stochastic gradient descent update for CRF

- For a training example $(\mathbf{x}_i, \mathbf{y}_i)$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha_t (\phi(\mathbf{x}_i, \mathbf{y}_i) - E_{\mathbf{y}}[\phi(\mathbf{x}_i, \mathbf{y})])$$

Structured perceptron

- For a training example $(\mathbf{x}_i, \mathbf{y}_i)$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha_t (\phi(\mathbf{x}_i, \mathbf{y}_i) - \phi(\mathbf{x}_i, \hat{\mathbf{y}}))$$

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$$

Expectation vs max



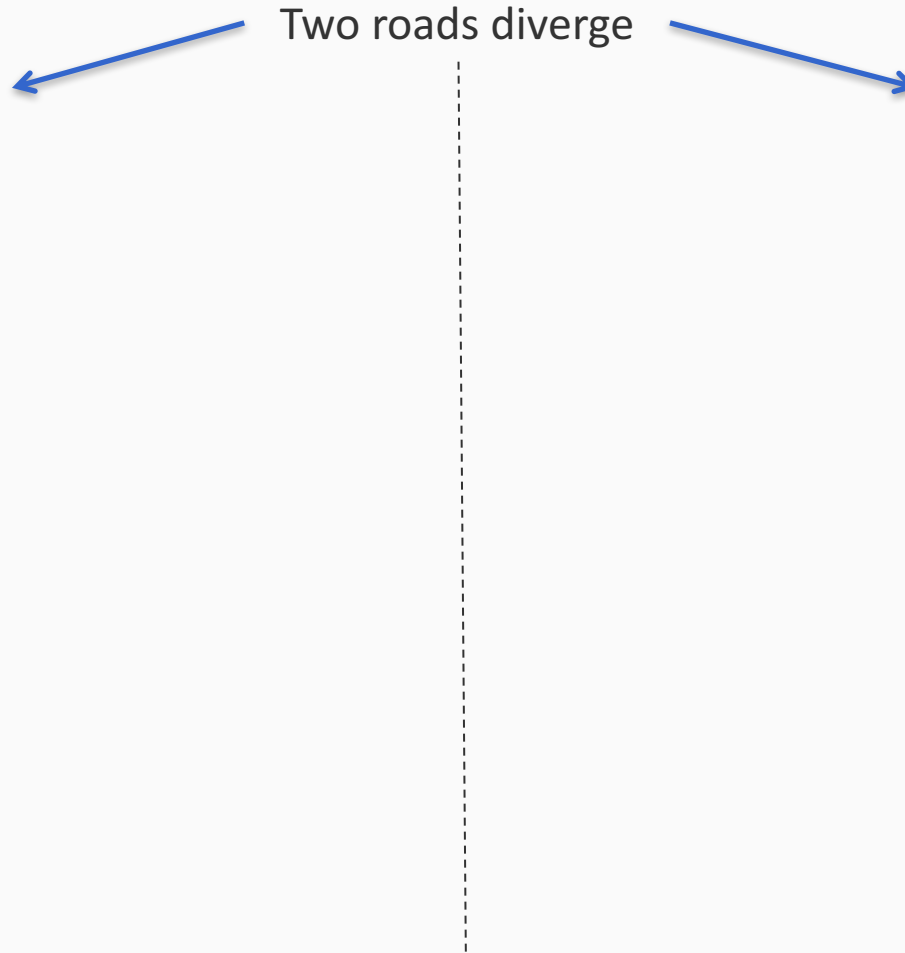
Caveat: Adding regularization will change the CRF update, averaging changes the perceptron update

The lay of the land

HMM: A generative model, assigns probabilities to sequences

The lay of the land

HMM: A generative model, assigns probabilities to sequences



The lay of the land

HMM: A generative model, assigns probabilities to sequences


Two roads diverge

- Hidden Markov Models are actually just **linear** classifiers
- Don't really care whether we are predicting probabilities. We are assigning scores to a full output for a given input (like multiclass)
- Generalize algorithms for linear classifiers. Sophisticated models that can use **arbitrary features**
- Structured Perceptron
Structured SVM

The lay of the land

HMM: A generative model, assigns probabilities to sequences

Two roads diverge



- Hidden Markov Models are actually just **linear** classifiers
- Don't really care whether we are predicting probabilities. We are assigning scores to a full output for a given input (like multiclass)
- Generalize algorithms for linear classifiers. Sophisticated models that can use **arbitrary features**
- Structured Perceptron
Structured SVM

- Model probabilities via exponential functions. Gives us the **log-linear** representation
- Log-probabilities for sequences for a given input
- Learn by maximizing likelihood. Sophisticated models that can use **arbitrary features**
- Conditional Random field

The lay of the land

HMM: A generative model, assigns probabilities to sequences

Two roads diverge

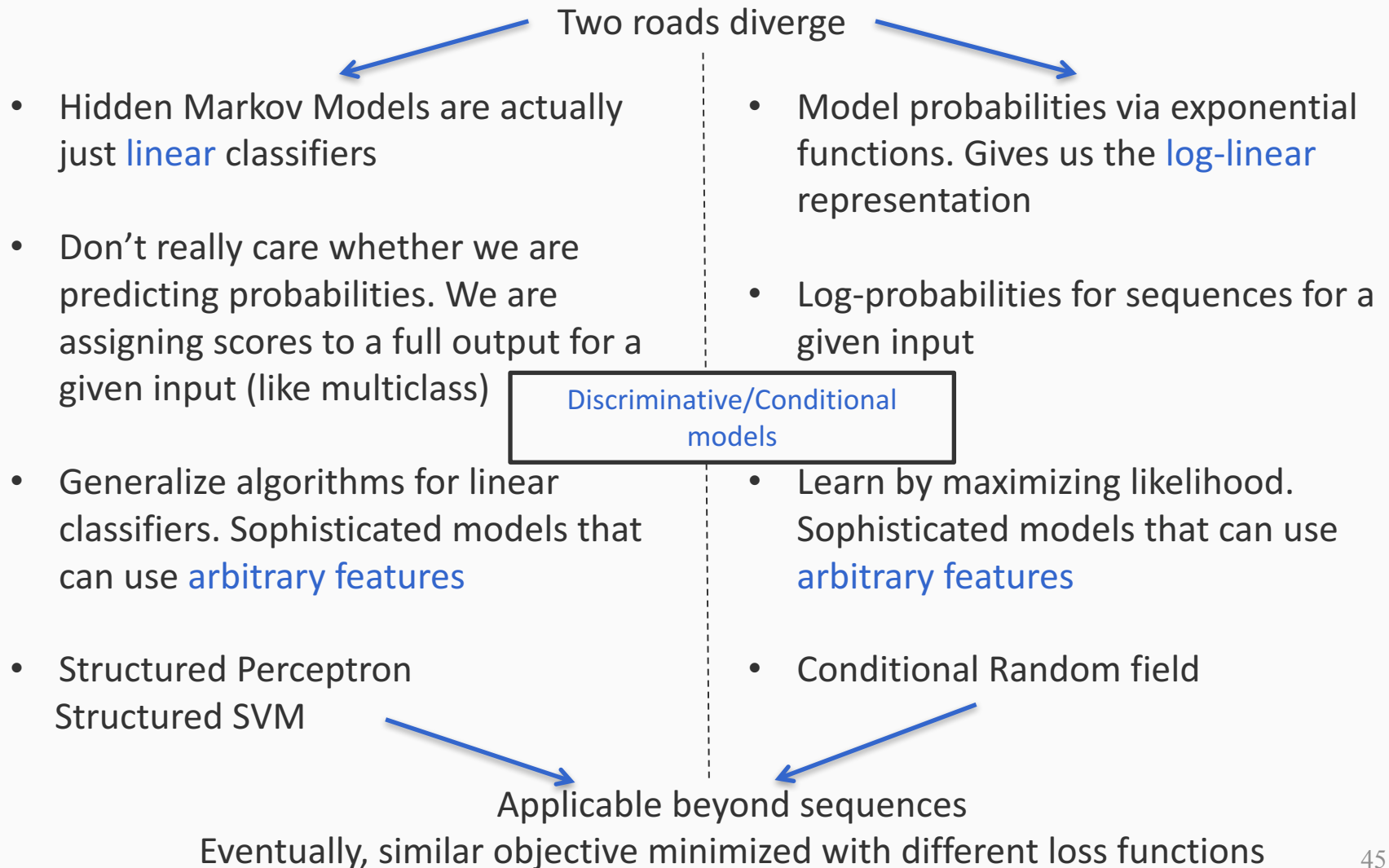
- Hidden Markov Models are actually just **linear** classifiers
- Don't really care whether we are predicting probabilities. We are assigning scores to a full output for a given input (like multiclass)
- Generalize algorithms for linear classifiers. Sophisticated models that can use **arbitrary features**
- Structured Perceptron
Structured SVM

- Model probabilities via exponential functions. Gives us the **log-linear** representation
- Log-probabilities for sequences for a given input
- Learn by maximizing likelihood. Sophisticated models that can use **arbitrary features**
- Conditional Random field

Discriminative/Conditional models

The lay of the land

HMM: A generative model, assigns probabilities to sequences



The lay of the land

HMM: A generative model, assigns probabilities to sequences

Two roads diverge

- Hidden Markov Models are actually just **linear** classifiers
- Don't really care whether we are predicting probabilities. We are assigning scores to a full output for a given input (like multiclass)
- Generalize algorithms for linear classifiers. Sophisticated models that can use **arbitrary features**
- Structured Perceptron
Structured SVM

- Model probabilities via exponential functions. Gives us the **log-linear** representation
- Log-probabilities for sequences for a given input
- Learn by maximizing likelihood. Sophisticated models that can use **arbitrary features**
- Conditional Random field

Discriminative/Conditional models

Applicable beyond sequences

Eventually, similar objective minimized with different loss functions

Coming soon...

Sequence models: Summary

- Goal: Predict an output sequence given input sequence
- Hidden Markov Model
- Inference
 - Predict via Viterbi algorithm
- Conditional models/discriminative models
 - **Local approaches** (no inference during training)
 - MEMM, conditional Markov model
 - **Global approaches** (inference during training)
 - CRF, structured perceptron
- To think
 - What are the *parts* in a sequence model?
 - How is each model scoring these parts?

Prediction is not always tractable for general structures

Same dichotomy for more general structures