

CSE 517

Natural Language Processing

Winter2015

Feature Rich Models

Yejin Choi - University of Washington

[Slides from Jason Eisner, Dan Klein, Luke Zettlemoyer]

Feature Rich Models

- **Throw anything you want into the stew**
- Add a bonus for this, a penalty for that, etc.



"11,001 New Features for Statistical Machine Translation",
(D. Chiang, K. Knight, and W. Wang), NAACL, 2009. Best
Paper Award.

Probabilistic Models

(Unstructured) categorization:

- Naïve Bayes

Structured prediction:

- HMMs
- PCFG Models
- IBM Models

Feature-rich / (Log)-linear Models

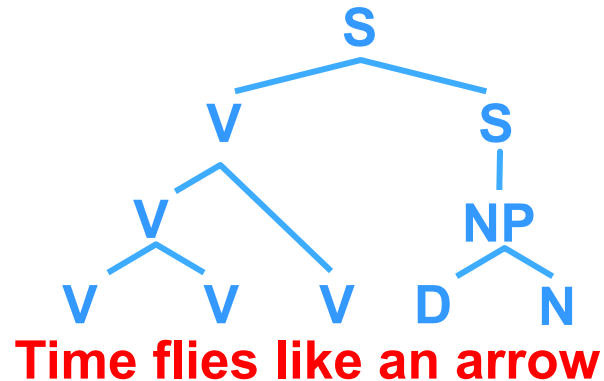
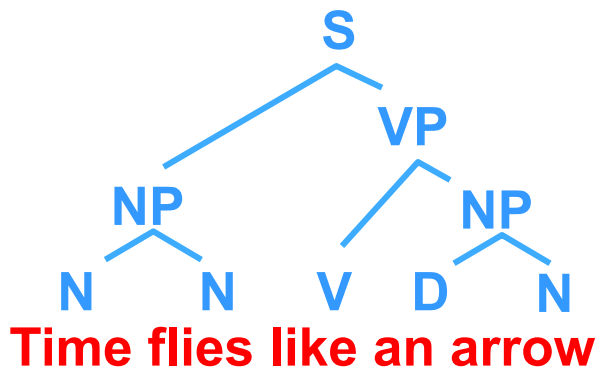
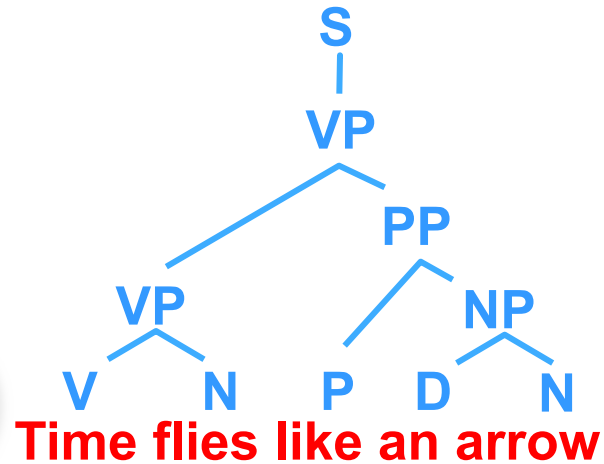
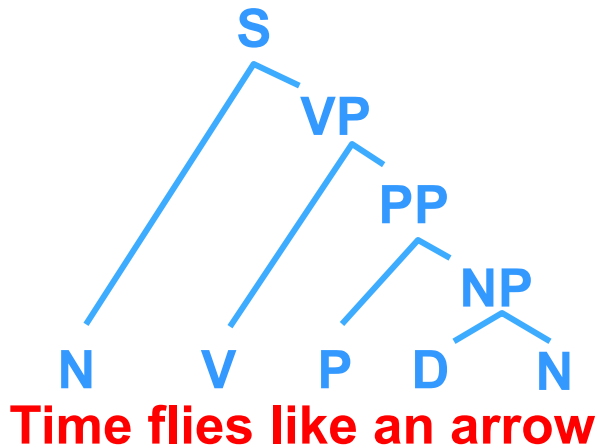
(Unstructured) categorization:

- Perceptron
- Maximum Entropy

Structured prediction:

- Perceptron for Structured Prediction
- MEMM (Maximum Entropy Markov Model)
- CRF (Conditional Random Fields)

Structured Prediction with Perceptrons and CRFs



...

$p(\text{category} \mid \text{message})$

goodmail

Reply today to claim your ...

spam

Reply today to claim your ...

goodmail

Wanna get pizza tonight?

spam

Wanna get pizza tonight?

goodmail

Thx; consider enlarging the ...

spam

Thx; consider enlarging the ...

goodmail

Enlarge your hidden ...

spam

Enlarge your hidden ...

p(RHS | LHS)

S → NP VP

S → N VP

S → NP[+wh] V S/V/NP

S → VP NP

S → Det N

S → PP P

...

p(RHS | LHS)

S → NP VP

S → N VP

S → NP[+wh] V S/V/NP

S → VP NP

S → Det N

S → PP P

...

NP → NP VP

NP → N VP

NP → NP CP/NP

NP → VP NP

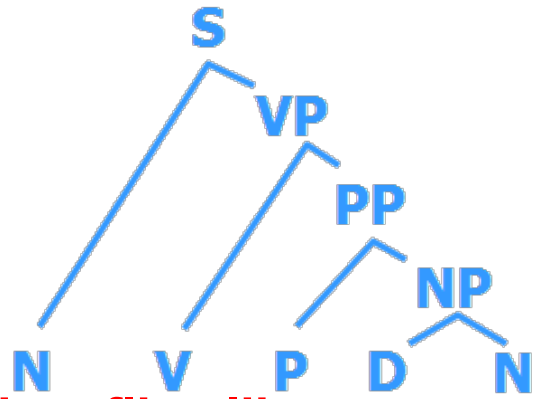
NP → Det N

NP → NP PP

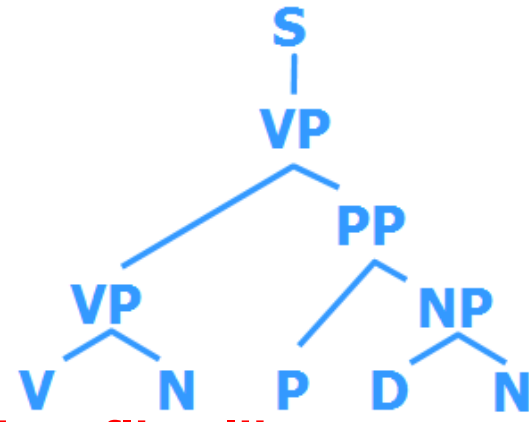
...

...

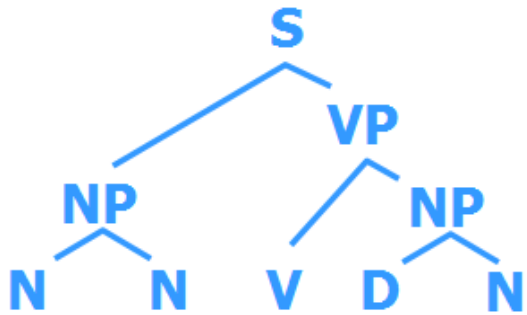
$p(\text{parse} \mid \text{sentence})$



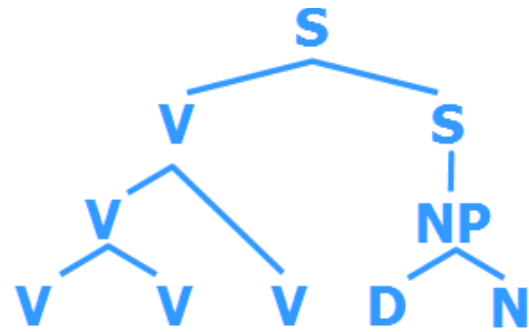
Time flies like an arrow



Time flies like an arrow



Time flies like an arrow



Time flies like an arrow

...

$p(\text{tag sequence} \mid \text{word sequence})$

N V P D N
Time flies like an arrow

V N P D N
Time flies like an arrow

...

N N V D N
Time flies like an arrow

V V V D N
Time flies like an arrow

Today's general problem

- Given some **input x**
 - Consider a set of **candidate outputs y**
 - **Classifications** for x (small number: often just 2)
 - **Taggings** of x (exponentially many)
 - **Parses** of x (exponential, even infinite)
 - **Translations** of x (exponential, even infinite)
 - ...
 - Want to find the "best" y , **given x**
- Structured prediction

Scoring by Linear Models

- Given some **input** x
- Consider a set of **candidate outputs** y
- Define a scoring function $\text{score}(x, y)$

Linear function: A sum of feature weights (you pick the features!)

Weight of feature k
(learned or set by hand)

$$\text{score}(x, y) = \sum_k \theta_k f_k(x, y)$$

Ranges over all features,
e.g., $k=5$ (numbered features)
or $k=\text{"see Det Noun"}$ (named features)

Whether (x, y) has feature k (0 or 1)
Or how many times it fires (≥ 0)
Or how strongly it fires (real #)

- Choose y that maximizes $\text{score}(x, y)$

Scoring by Linear Models

- Given some **input** x
- Consider a set of **candidate outputs** y
- Define a scoring function $\text{score}(x, y)$

Linear function: A sum of feature weights (you pick the features!)

(learned or set by hand)

$$\text{score}(x, y) = \vec{\theta} \cdot \vec{f}(x, y)$$

This linear decision rule is sometimes called a “perceptron.”
It’s a “structured perceptron” if it does structured prediction
(number of y candidates is unbounded, e.g., grows with $|x|$).

- Choose y that maximizes $\text{score}(x, y)$

Probabilistic Models

(Unstructured) categorization:

- Naïve Bayes

Structured prediction:

- HMMs
- PCFG Models
- IBM Models

Feature-rich / (Log)-linear Models

(Unstructured) categorization:

Perceptron

- Maximum Entropy

Structured prediction:

- Perceptron for Structured Prediction
- MEMM (Maximum Entropy Markov Model)
- CRF (Conditional Random Fields)

Perceptron Training Algorithm

- initialize θ (usually to the zero vector)
- repeat:
 - Pick a training example (x, y)
 - Model predicts y^* that maximizes $\text{score}(x, y^*)$
 - Update weights by a step of size $\varepsilon > 0$:
$$\theta = \theta + \varepsilon \cdot (f(x, y) - f(x, y^*))$$

If model prediction was correct ($y=y^*$), θ doesn't change.
So once model predicts all training examples correctly, stop.

If some θ can do the job, this eventually happens!

(If not, θ will oscillate, but the average θ from all steps will settle down. So return that eventual average.)

Perceptron Training Algorithm

- initialize θ (usually to the zero vector)
- repeat:
 - Pick a training example (\mathbf{x}, y)
 - Model predicts y^* that maximizes $\text{score}(\mathbf{x}, y^*)$
 - Update weights by a step of size $\varepsilon > 0$:
$$\theta = \theta + \varepsilon \cdot (f(\mathbf{x}, y) - f(\mathbf{x}, y^*))$$

If model prediction was wrong ($y \neq y^*$), then we must have $\text{score}(\mathbf{x}, y) \leq \text{score}(\mathbf{x}, y^*)$ instead of $>$ as we want.

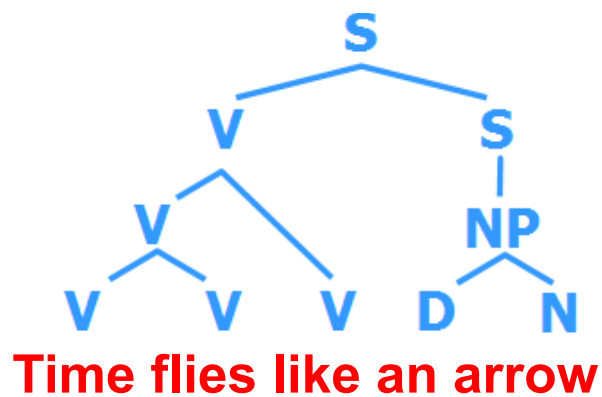
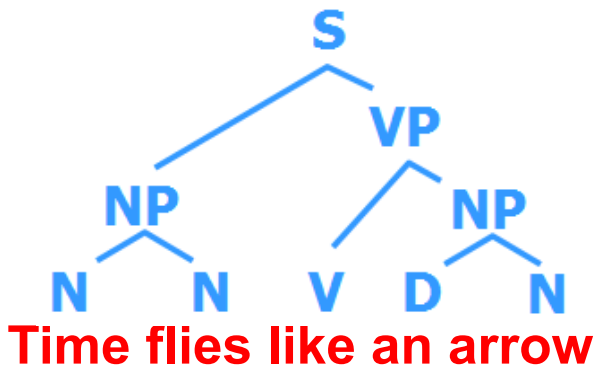
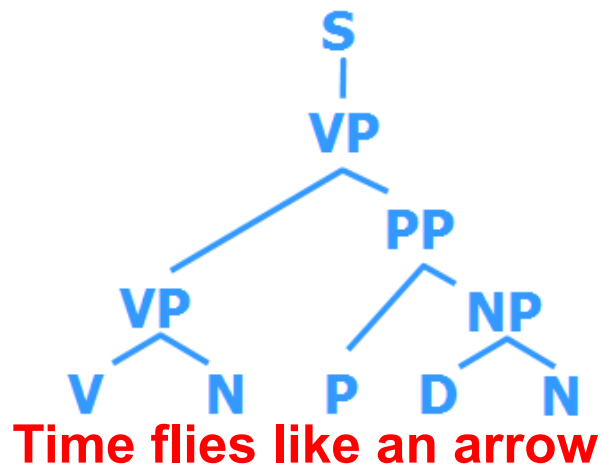
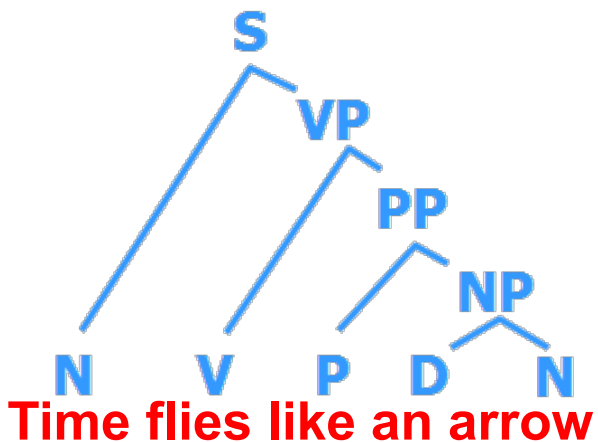
Equivalently, $\theta \cdot f(\mathbf{x}, y) \leq \theta \cdot f(\mathbf{x}, y^*)$

Equivalently, $\theta \cdot (f(\mathbf{x}, y) - f(\mathbf{x}, y^*)) \leq 0$ but we want it positive.

Our update increases it (by $\varepsilon \cdot \|f(\mathbf{x}, y) - f(\mathbf{x}, y^*)\|^2 \geq 0$)

~~p(parse | sentence)~~

score(sentence, parse)



...

Nuthin' but adding weights

- n-grams: ... + $\log p(w_7 \mid w_5, w_6)$ + $\log p(w_8 \mid w_6, w_7)$ + ...
- PCFG: $\log p(\text{NP VP} \mid S)$ + $\log p(\text{Papa} \mid \text{NP})$ + $\log p(\text{VP PP} \mid \text{VP})$...
- HMM tagging: ... + $\log p(t_7 \mid t_5, t_6)$ + $\log p(w_7 \mid t_7)$ + ...
- Noisy channel: $[\log p(\text{source})]$ + $[\log p(\text{data} \mid \text{source})]$
- Cascade of composed FSTs:
 $[\log p(A)] + [\log p(B \mid A)] + [\log p(C \mid B)] + \dots$
- Naïve Bayes:
 $\log p(\text{Class}) + \log p(\text{feature1} \mid \text{Class}) + \log p(\text{feature2} \mid \text{Class})$...

What if our weights were arbitrary real numbers?

Change $\log p(\text{this} \mid \text{that})$ to $\theta(\text{this} ; \text{that})$

- **n-grams:** ... + $\log p(w_7 \mid w_5, w_6)$ + $\log p(w_8 \mid w_6, w_7)$ + ...
- **PCFG:** $\log p(\text{NP VP} \mid S)$ + $\log p(\text{Papa} \mid \text{NP})$ + $\log p(\text{VP PP} \mid \text{VP})$...
- **HMM tagging:** ... + $\log p(t_7 \mid t_5, t_6)$ + $\log p(w_7 \mid t_7)$ + ...
- **Noisy channel:** $[\log p(\text{source})]$ + $[\log p(\text{data} \mid \text{source})]$
- **Cascade of FSTs:**
 $[\log p(A)] + [\log p(B \mid A)] + [\log p(C \mid B)] + \dots$
- **Naïve Bayes:**
 $\log p(\text{Class}) + \log p(\text{feature1} \mid \text{Class}) + \log p(\text{feature2} \mid \text{Class}) \dots$

What if our weights were arbitrary real numbers?

Change $\log p(\text{this} \mid \text{that})$ to $\theta(\text{this} ; \text{that})$

- n-grams: ... + $\theta(w7 ; w5, w6)$ + $\theta(w8 ; w6, w7)$ + ...
- PCFG: $\theta(\text{NP VP} ; \text{S})$ + $\theta(\text{Papa} ; \text{NP})$ + $\theta(\text{VP PP} ; \text{VP})$...
- HMM tagging: ... + $\theta(t7 ; t5, t6)$ + $\theta(w7 ; t7)$ + ...
- Noisy channel: $[\theta(\text{source})]$ + $[\theta(\text{data} ; \text{source})]$
- Cascade of FSTs:
 $[\theta(A)] + [\theta(B ; A)] + [\theta(C ; B)] + \dots$
- Naïve Bayes:
 $\theta(\text{Class}) + \theta(\text{feature1} ; \text{Class}) + \theta(\text{feature2} ; \text{Class}) \dots$

In practice, θ is a hash table

Maps from feature name (a string or object) to feature weight (a float)

e.g., $\theta(\text{NP VP} ; \text{S})$ = weight of the $\text{S} \rightarrow \text{NP VP}$ rule, say -0.1 or +1.3

What if our weights were arbitrary real numbers?

Change $\log p(\text{this} \mid \text{that})$ to $\theta(\text{this} ; \text{that})$

- n-grams: ... + $\theta(w_5 w_6 w_7)$ + $\theta(w_6 w_7 w_8)$ + ...

WCFG

- ~~PCFG~~: $\theta(S \rightarrow NP VP)$ + $\theta(NP \rightarrow \text{Papa})$ + $\theta(VP \rightarrow VP PP)$...

- HMM tagging: ... + $\theta(t_5 t_6 t_7)$ + $\theta(t_7 \rightarrow w_7)$ + ...

- Noisy channel: [$\theta(\text{source})$] + [$\theta(\text{source}, \text{data})$]

- Cascade of FSTs:

$$[\theta(A)] + [\theta(A, B)] + [\theta(B, C)] + \dots$$

- ~~Naïve Bayes~~: (multi-class) logistic regression
 $\theta(\text{Class})$ + $\theta(\text{Class}, \text{feature 1})$ + $\theta(\text{Class}, \text{feature 2})$...

Finding the best y given x

- At **both** training & test time, given input x , perceptron picks y that maximizes $\text{score}(x,y)$

$$\text{score}(x, y) = \sum_k \theta_k f_k(x, y)$$

- How do we find $\text{argmax}_y \text{score}(x,y)$?
 - Easy when only a few candidates y (e.g., text classification)
 - Just try each y in turn.
 - Harder for structured prediction: **but you now know how!**
 - Find the **best string, path, or tree ...**
 - Viterbi for HMM, CKY for trees, stack decoding for MT
 - Dynamic programming if possible

Why would we switch from probabilities to scores?

1. “Discriminative” training (e.g., perceptron) might work better.

- It tries to optimize weights to actually predict the right y for each x .
- More important than maximizing $\log p(x,y) = \log p(y|x) + \log p(x)$, as we’ve been doing in HMMs and PCFGs.
- Satisfied once the right y wins. The example puts no more pressure on the weights to raise $\log p(y|x)$. And never pressures us to raise $\log p(x)$.

2. Having more freedom in the weights might help?

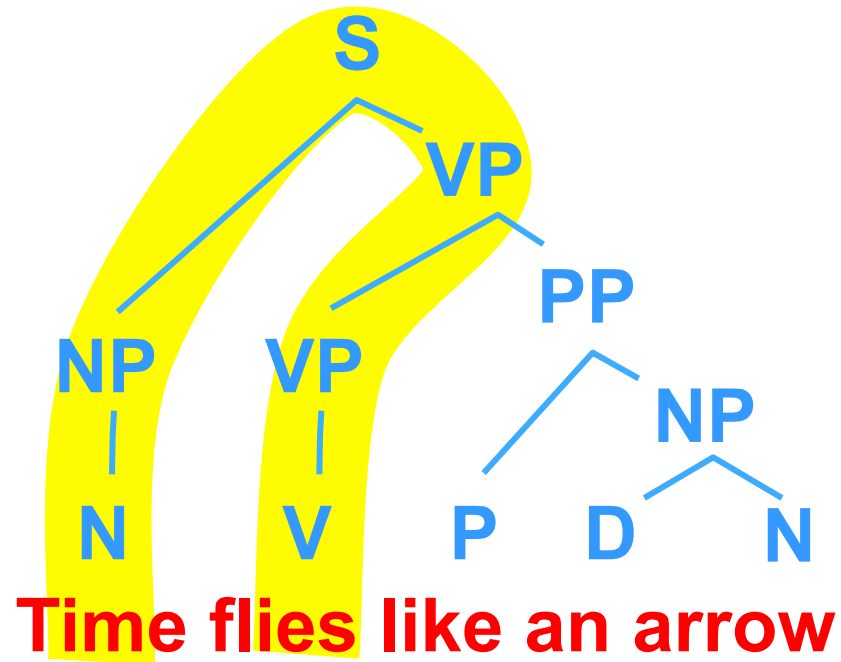
- Now weights can be positive or negative.
- Exponentiated weights no longer have to sum to 1.
- But turns out new θ vectors can’t do more than the old restricted ones.
 - Roughly, for every WCFG there’s an equivalent PCFG.
 - Though it’s true a regularizer might favor one of the new ones.

3. We can throw lots more features into the stewpot.

- Allows model to capture more of the useful predictive patterns!
- **So, what features can we throw in efficiently?**

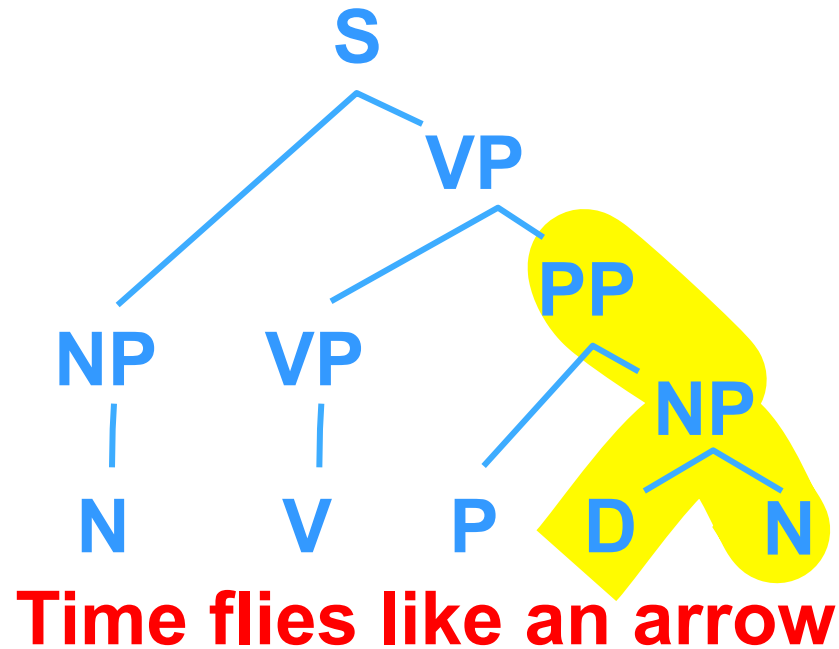


Cross-rule substructures



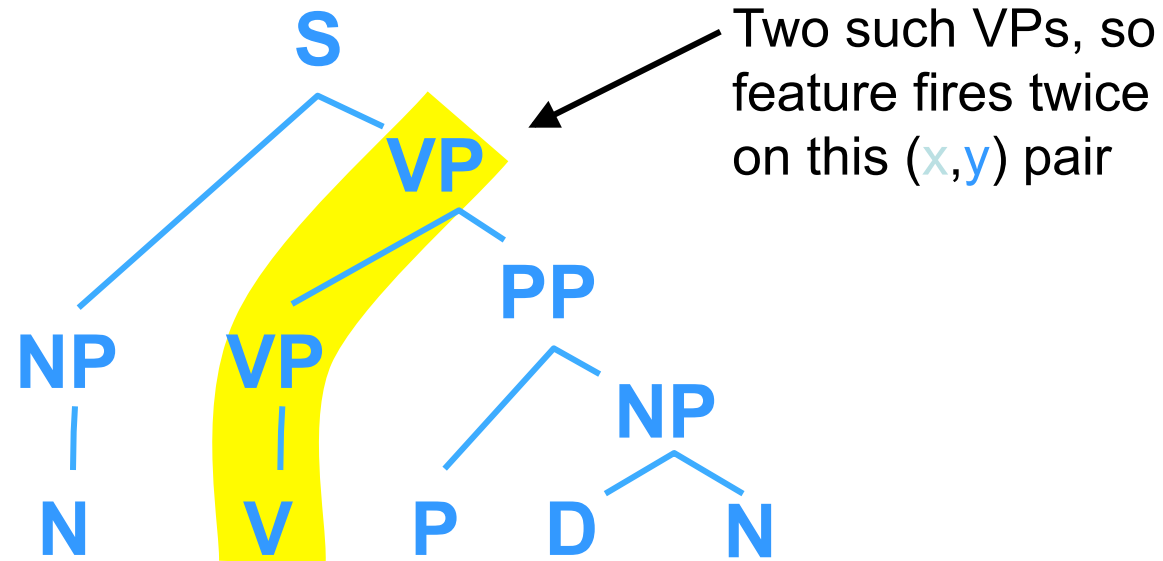
- Count of “flies” as a verb with subject “time”

Cross-rule substructures



- Count of “flies” as a verb with subject “time”
- Count of NP → D N when the NP is the object of a preposition

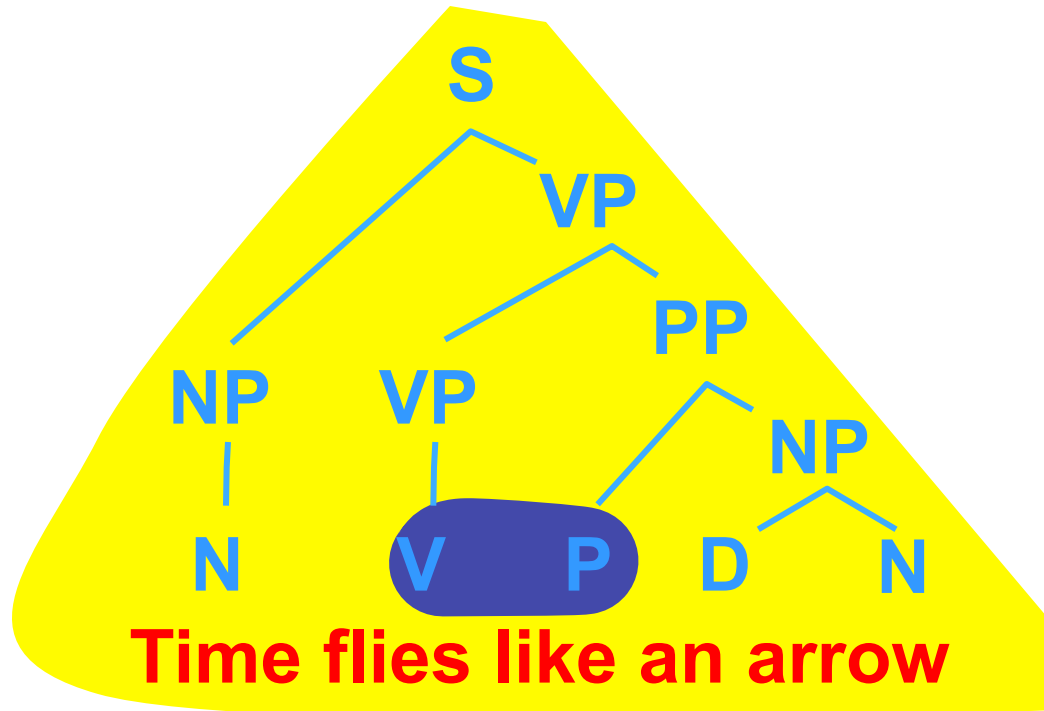
Cross-rule substructures



Time flies like an arrow

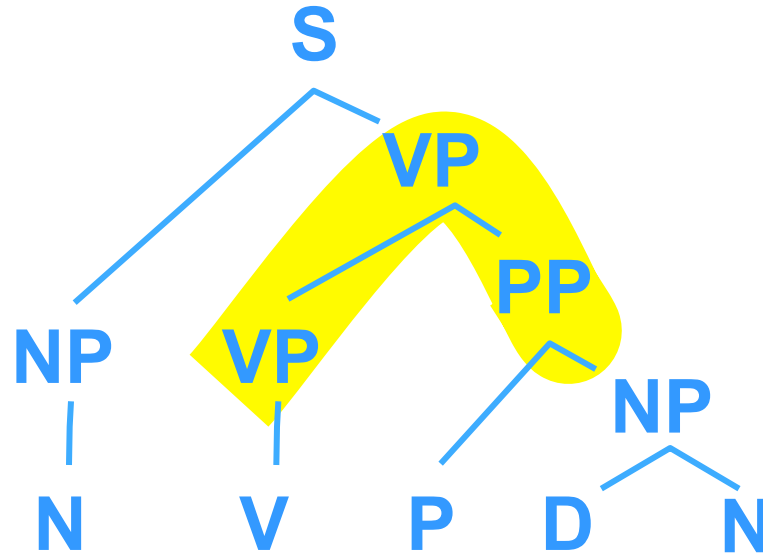
- Count of “flies” as a verb with subject “time”
- Count of NP → D N when the NP is the object of a preposition
- Count of VPs that contain a V

Global features



- Count of “NP and NP” when the two NPs have very different size or structure [this feature has weight < 0]
- The number of PPs is even
- The depth of the tree is prime 😊
- Count of the tag bigram V P in the preterminal seq

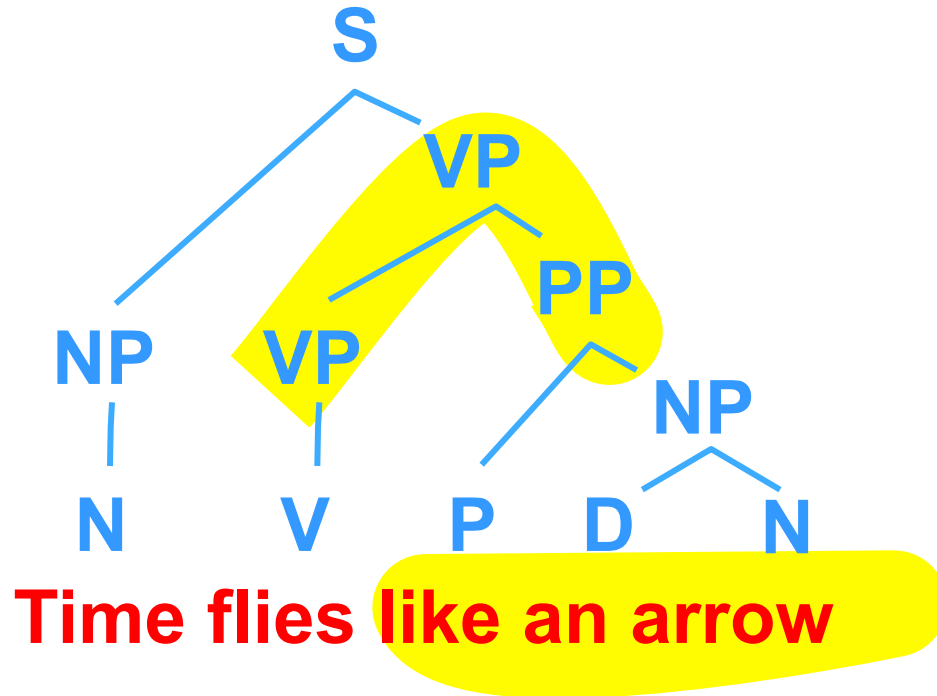
Context-specific features



Time flies like an arrow

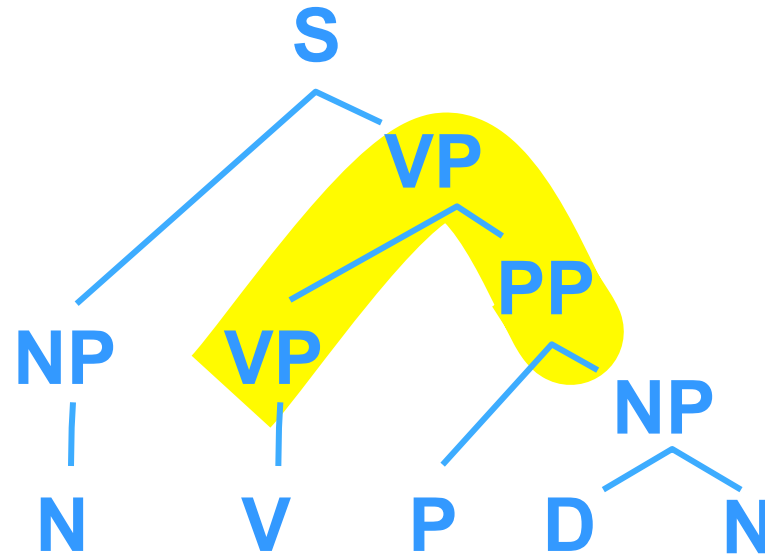
- Count of VP → VP PP whose first word is “flies”

Context-specific features



- Count of VP → VP PP whose first word is “flies”
- Count of VP → VP PP whose right child has width 3

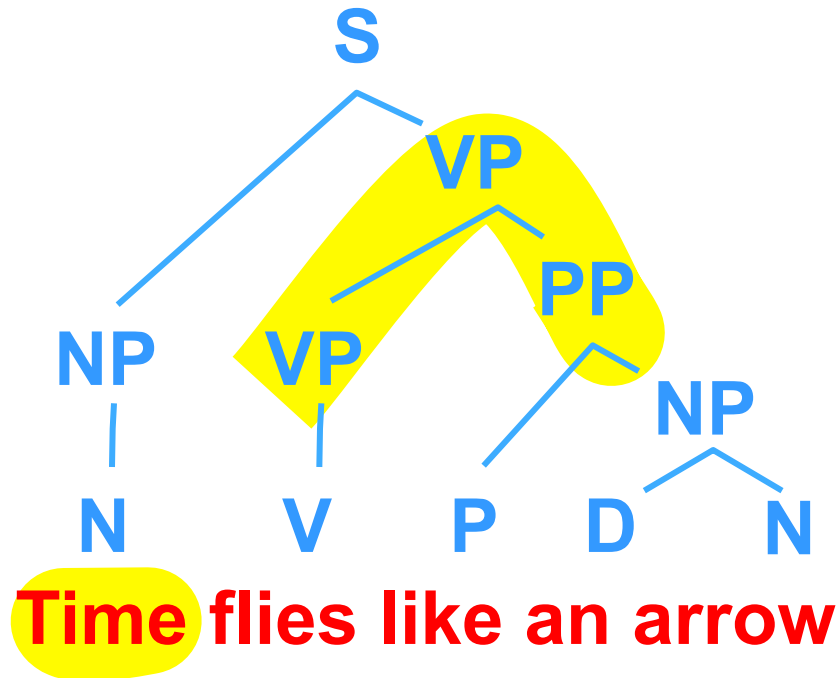
Context-specific features



Time flies like an arrow 5

- Count of VP → VP PP whose first word is “flies”
- Count of VP → VP PP whose right child has width 3
- Count of VP → VP PP at the end of the input

Context-specific features



- Count of VP → VP PP whose first word is “flies”
- Count of VP → VP PP whose right child has width 3
- Count of VP → VP PP at the end of the input
- Count of VP → VP PP right after a capitalized word

In the case of tagging ...

N V P D N
Time flies like an arrow

- Count of tag P as the tag for “like”
- Count of tag P
- Count of tag P in the middle third of the sentence
- Count of tag bigram V P
- Count of tag bigram V P followed by “an”
- Count of tag bigram V P where P is the tag for “like”
- Count of tag bigram V P where both words are lowercase

Overview: POS tagging Accuracies

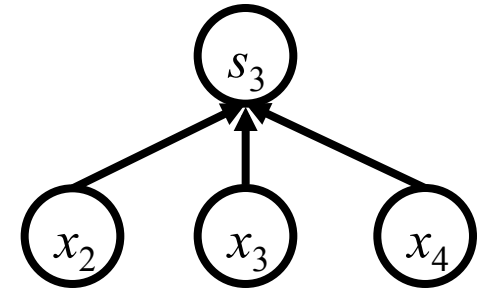
- Roadmap of (known / unknown) accuracies:
 - Most freq tag: ~90% / ~50%
 - Trigram HMM: ~95% / ~55%
 - TnT (HMM++): 96.2% / 86.0%

- What if feature-rich classifier that predicts each POS tag one at a time?

- Upper bound: ~98%

What about better features?

- Choose the most common tag
 - 90.3% with a bad unknown word model
 - 93.7% with a good one
- What about looking at a word and its environment, but no sequence information?
 - Add in previous / next word
 - Previous / next word shapes
 - Occurrence pattern features
 - Crude entity detection
 - Phrasal verb in sentence?
 - Conjunctions of these things
- Uses lots of features: > 200K



the ___
X ___ X
[X: x X occurs]
___ (Inc.|Co.)
put ___

Probabilistic Models

(Unstructured) categorization:

- Naïve Bayes

Structured prediction:

- HMMs
- PCFG Models
- IBM Models

Feature-rich / (Log)-linear Models

(Unstructured) categorization:

- 
- Perceptron
 - Maximum Entropy

Structured prediction:

- Perceptron for Structured Prediction
- MEMM (Maximum Entropy Markov Model)
- CRF (Conditional Random Fields)

Maximum Entropy (MaxEnt) Models

- Also known as “Log-linear” Models (*linear if you take log*)

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}^\top \mathbf{f}(\mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{y}'))}$$

- The feature vector representation may include redundant and overlapping features

Training MaxEnt Models

- Maximizing the likelihood of the training data incidentally maximizes the entropy (hence “maximum entropy”)

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}^\top \mathbf{f}(\mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{y}'))}$$

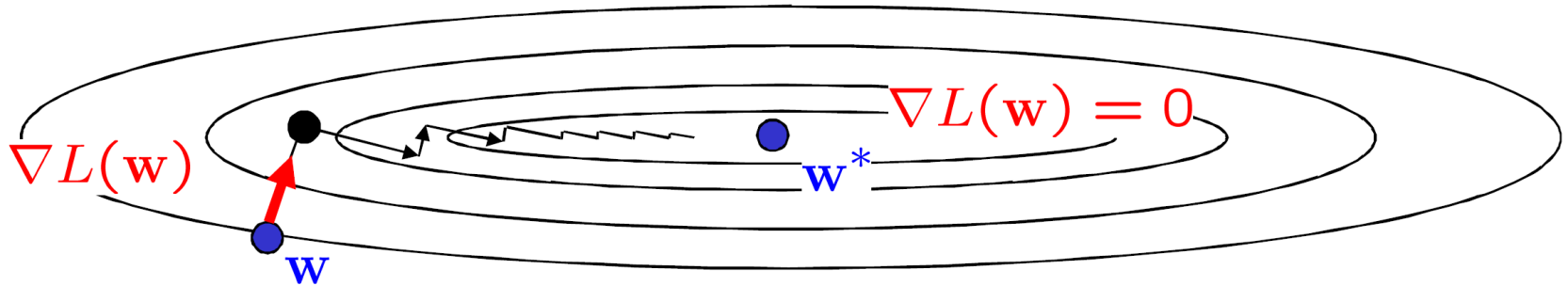
← Make positive
← Normalize

- Maximize the (log) conditional likelihood of training data

$$L(\mathbf{w}) = \log \prod_i P(\mathbf{y}^i | \mathbf{x}^i, \mathbf{w}) = \sum_i \log \left(\frac{\exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i))}{\sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}))} \right)$$
$$= \sum_i \left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

Convex Optimization for Training

$$L(\mathbf{w})$$



- The likelihood function is convex. (can get global optimum)
- Many optimization algorithms/software available.
 - Gradient ascent (descent), Conjugate Gradient, L-BFGS, etc
- All we need are:
 - (1) evaluate the function at current 'w'
 - (2) evaluate its derivative at current 'w'

Training MaxEnt Models

$$L(\mathbf{w}) = \sum_i \left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}_n} = \sum_i \left(\mathbf{f}_i(\mathbf{y}^i)_n - \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}_i) \mathbf{f}_i(\mathbf{y})_n \right)$$



Total count of feature n
in correct candidates



Expected count of
feature n in predicted
candidates

Training with Regularization

$$L(\mathbf{w}) = -k\|\mathbf{w}\|^2 + \sum_i \left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_y \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}_n} = -2k\mathbf{w}_n + \sum_i \left(\mathbf{f}_i(\mathbf{y}^i)_n - \sum_y P(\mathbf{y}|\mathbf{x}_i) \mathbf{f}_i(\mathbf{y})_n \right)$$


Big weights are bad



Total count of feature n
in correct candidates

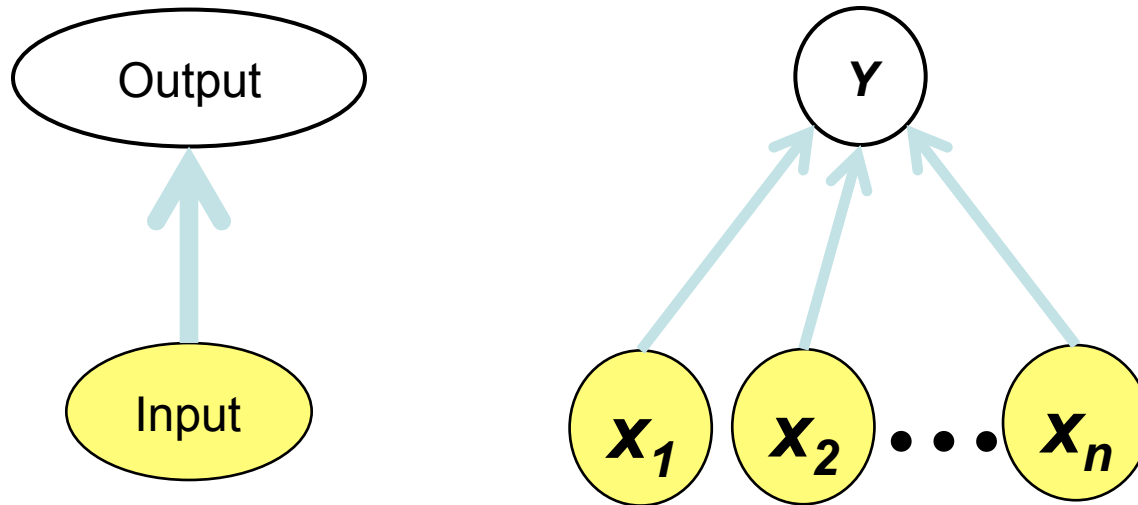


Expected count of
feature n in predicted
candidates



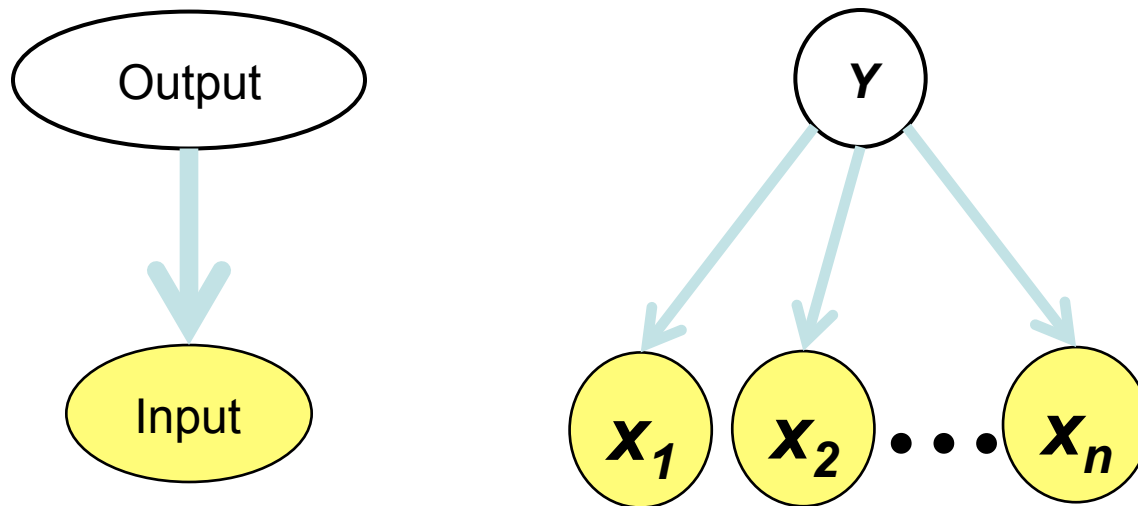
Graphical Representation of MaxEnt

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}^\top \mathbf{f}(\mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{y}'))}$$

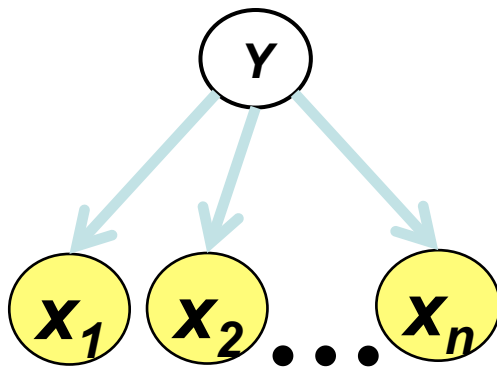


Graphical Representation of Naïve Bayes

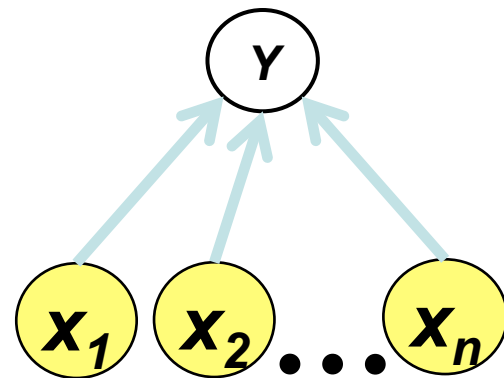
$$P(X | Y) = \prod_{j=1} P(x_j | Y)$$



Naïve Bayes



MaxEnt



Naïve Bayes Classifier

“**Generative**” models

- $p(\text{input} \mid \text{output})$
- For instance, for text categorization, $P(\text{words} \mid \text{category})$
- Unnecessary efforts on generating input

→ Independent assumption among input variables: Given the category, each word is generated independently from other words (too strong assumption in reality!)

→ Cannot incorporate arbitrary/redundant/overlapping features

Maximum Entropy Classifier

“**Discriminative**” models

- $p(\text{output} \mid \text{input})$
- For instance, for text categorization, $P(\text{category} \mid \text{words})$
- Focus directly on predicting the output

→ By conditioning on the entire input, we don't need to worry about the independent assumption among input variables

→ Can incorporate arbitrary features: redundant and overlapping features

Overview: POS tagging Accuracies

- Roadmap of (known / unknown) accuracies:
 - Most freq tag: ~90% / ~50%
 - Trigram HMM: ~95% / ~55%
 - TnT (HMM++): 96.2% / 86.0%
 - Maxent $P(s_i|x)$: 96.8% / 86.8%
- Q: What does this say about sequence models?
- Q: How do we add more features to our sequence models?
- Upper bound: ~98%

Probabilistic Models

(Unstructured) categorization:

- Naïve Bayes

Structured prediction:

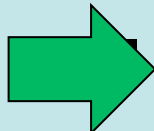
- HMMs
- PCFG Models
- IBM Models

Feature-rich / (Log)-linear Models

(Unstructured) categorization:

- 
- Perceptron
 - Maximum Entropy

Structured prediction:

- 
- Perceptron for Structured Prediction
 - MEMM (Maximum Entropy Markov Model)
 - CRF (Conditional Random Fields)

MEMM Taggers

- **One step up:** also condition on previous tags

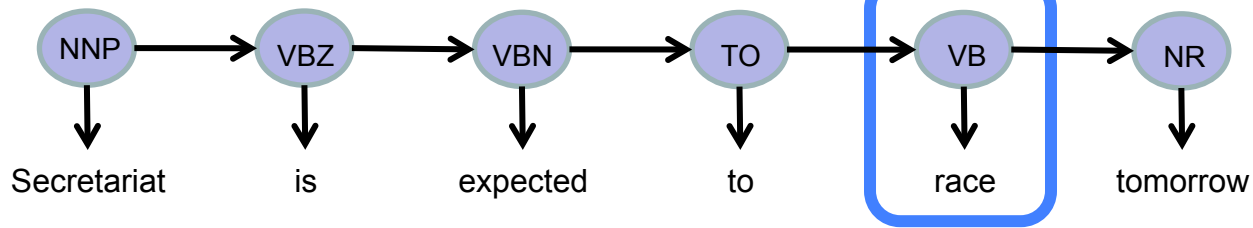
$$\begin{aligned} p(s_1 \dots s_m | x_1 \dots x_m) &= \prod_{i=1}^m p(s_i | s_1 \dots s_{i-1}, x_1 \dots x_m) \\ &= \prod_{i=1}^m p(s_i | s_{i-1}, x_1 \dots x_m) \end{aligned}$$

- Train up $p(s_i | s_{i-1}, x_1 \dots x_m)$ as a discrete log-linear (maxent) model, then use to score sequences

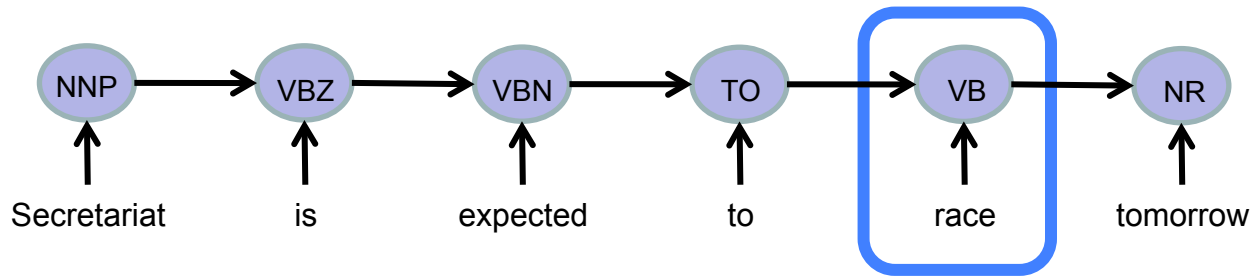
$$p(s_i | s_{i-1}, x_1 \dots x_m) = \frac{\exp(w \cdot \phi(x_1 \dots x_m, i, s_{i-1}, s_i))}{\sum_{s'} \exp(w \cdot \phi(x_1 \dots x_m, i, s_{i-1}, s'))}$$

- This is referred to as an MEMM tagger [Ratnaparkhi 96]
- Beam search effective! (Why?)
- What's the advantage of beam size 1?

HMM



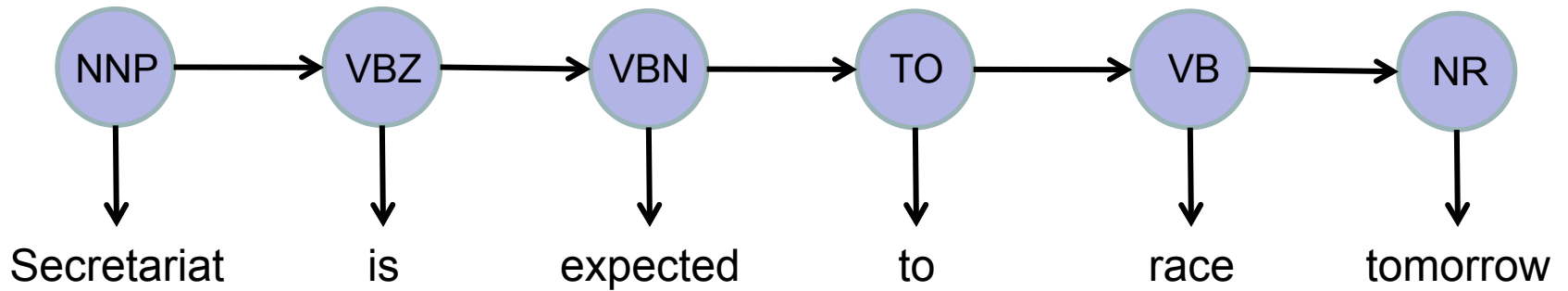
MEMM



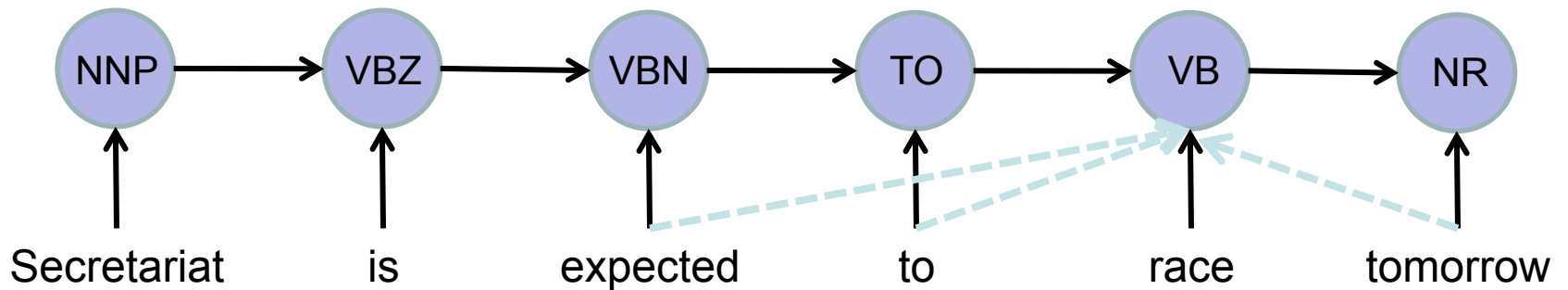
HMM	MEMM
<p>“Generative” models</p> <ul style="list-style-type: none"> → joint probability $p(\text{words, tags})$ → “generate” input (in addition to tags) → but we need to predict tags, not words! 	<p>“Discriminative” or “Conditional” models</p> <ul style="list-style-type: none"> → conditional probability $p(\text{tags} \text{words})$ → “condition” on input → Focusing only on predicting tags
<p>Probability of each slice = emission * transition = $p(\text{word}_i \text{tag}_i) * p(\text{tag}_i \text{tag}_{i-1}) =$</p> <p>→ Cannot incorporate long distance features</p>	<p>Probability of each slice = $p(\text{tag}_i \text{tag}_{i-1}, \text{word}_i)$ or $p(\text{tag}_i \text{tag}_{i-1}, \text{all words})$</p> <p>→ Can incorporate long distance features</p>

HMM v.s. MEMM

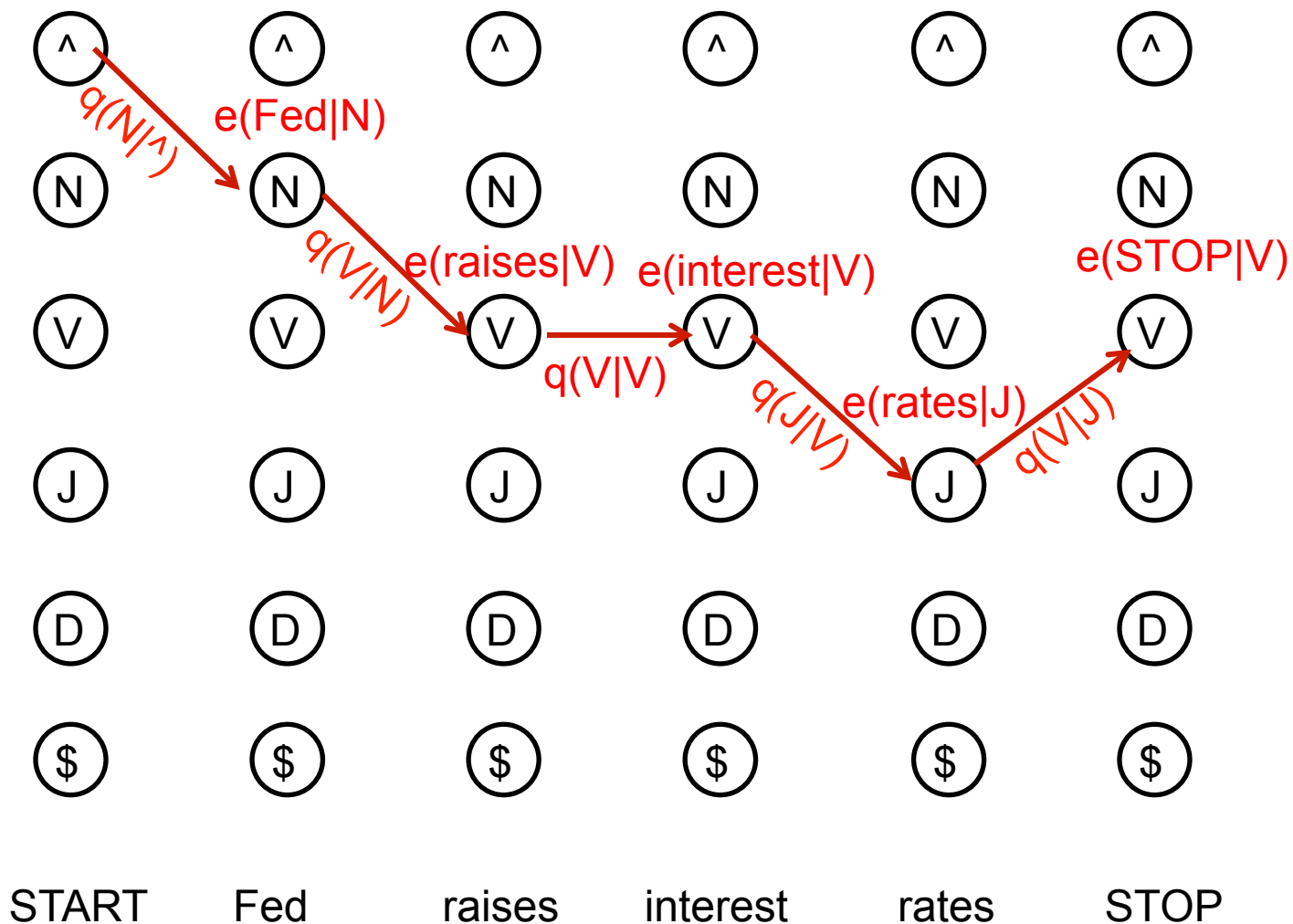
HMM



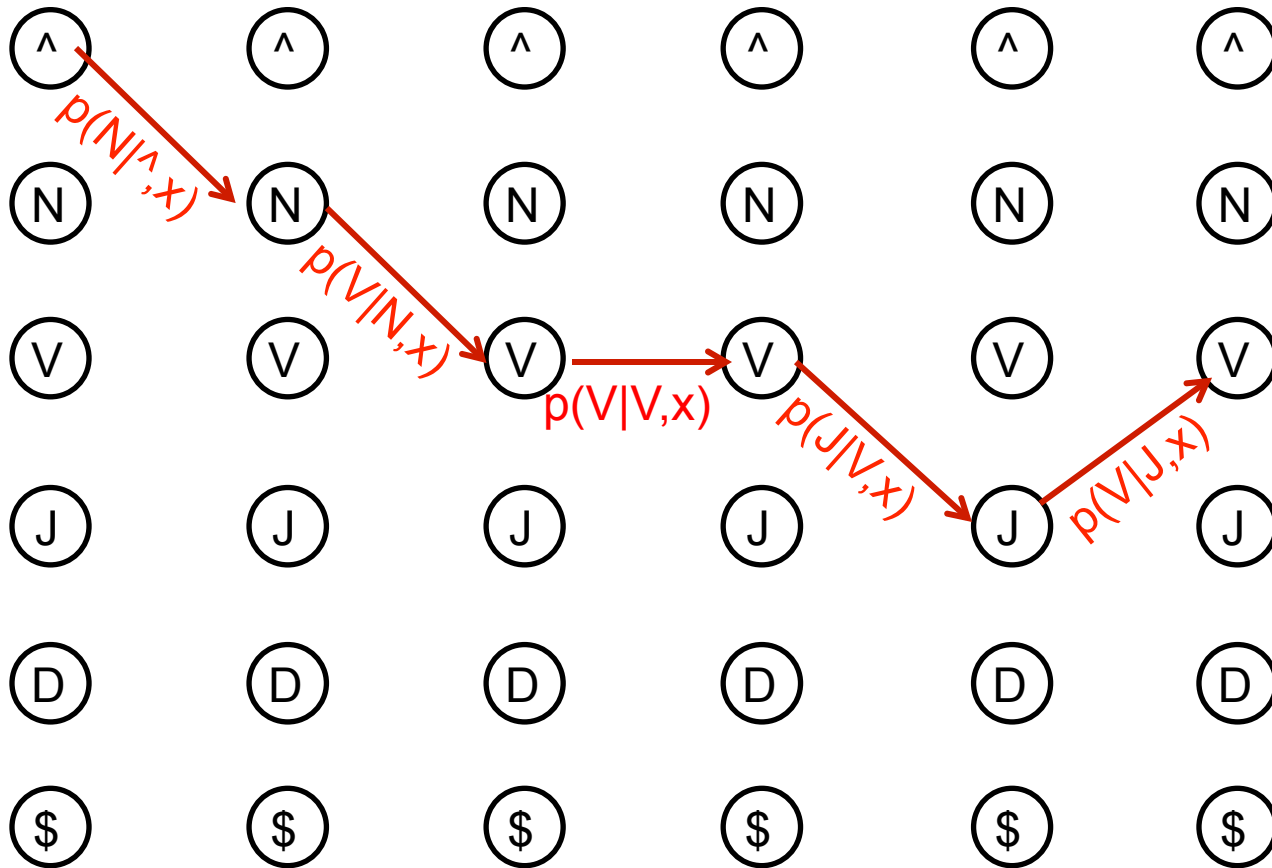
MEMM



The HMM State Lattice / Trellis (repeat slide)



The MEMM State Lattice / Trellis



x = START

Fed

raises

interest

rates

STOP

Decoding: $p(s_1 \dots s_m | x_1 \dots x_m) = \prod_{i=1}^m p(s_i | s_1 \dots s_{i-1}, x_1 \dots x_m)$

- **Decoding maxent taggers:**

- Just like decoding HMMs
- Viterbi, beam search, posterior decoding

- **Viterbi algorithm (HMMs):**

- Define $\pi(i, s_i)$ to be the max score of a sequence of length i ending in tag s_i

$$\pi(i, s_i) = \max_{s_{i-1}} e(x_i | s_i) q(s_i | s_{i-1}) \pi(i-1, s_{i-1})$$

- **Viterbi algorithm (Maxent):**

- Can use same algorithm for MEMMs, just need to redefine $\pi(i, s_i)$!

$$\pi(i, s_i) = \max_{s_{i-1}} p(s_i | s_{i-1}, x_1 \dots x_m) \pi(i-1, s_{i-1})$$

Overview: Accuracies

- Roadmap of (known / unknown) accuracies:
 - Most freq tag: ~90% / ~50%
 - Trigram HMM: ~95% / ~55%
 - TnT (HMM++): 96.2% / 86.0%
 - Maxent $P(s_i|x)$: 96.8% / 86.8%
 - MEMM tagger: 96.9% / 86.9%

- Upper bound: ~98%

Global Discriminative Taggers

- **Newer, higher-powered discriminative sequence models**
 - CRFs (also perceptrons, M3Ns)
 - Do not decompose training into independent local regions
 - Can be deathly slow to train – require repeated inference on training set
- **Differences can vary in importance, depending on task**
- **However:** one issue worth knowing about in local models
 - “**Label bias**” and other **explaining away effects**
 - MEMM taggers’ local scores can be near one without having both good “transitions” and “emissions”
 - This means that often evidence doesn’t flow properly
 - Why isn’t this a big deal for POS tagging?
 - Also: in decoding, condition on predicted, not gold, histories

Probabilistic Models

(Unstructured) categorization:

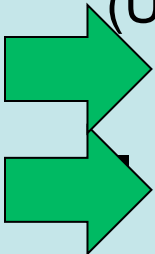
- Naïve Bayes

Structured prediction:


- HMMs
- PCFG Models
- IBM Models

Feature-rich / (Log)-linear Models

(Unstructured) categorization:

- 
- Perceptron
 - Maximum Entropy

Structured prediction:

- 
- Perceptron for Structured Prediction
 - MEMM (Maximum Entropy Markov Model)
 - CRF (Conditional Random Fields)

Linear Models: Perceptron

- The perceptron algorithm
 - Iteratively processes the training set, reacting to training errors
 - Can be thought of as trying to drive down training error
- The (online) perceptron algorithm:

- Start with zero weights
- Visit training instances (x_i, y_i) one by one
 - Make a prediction

$$y^* = \arg \max_y w \cdot \phi(x_i, y)$$

Sentence: $x = x_1 \dots x_m$

Tag Sequence:
 $y = s_1 \dots s_m$

- If correct ($y^* = y_i$): no change, goto next example!
- If wrong: adjust weights

$$w = w + \phi(x_i, y_i) - \phi(x_i, y^*)$$

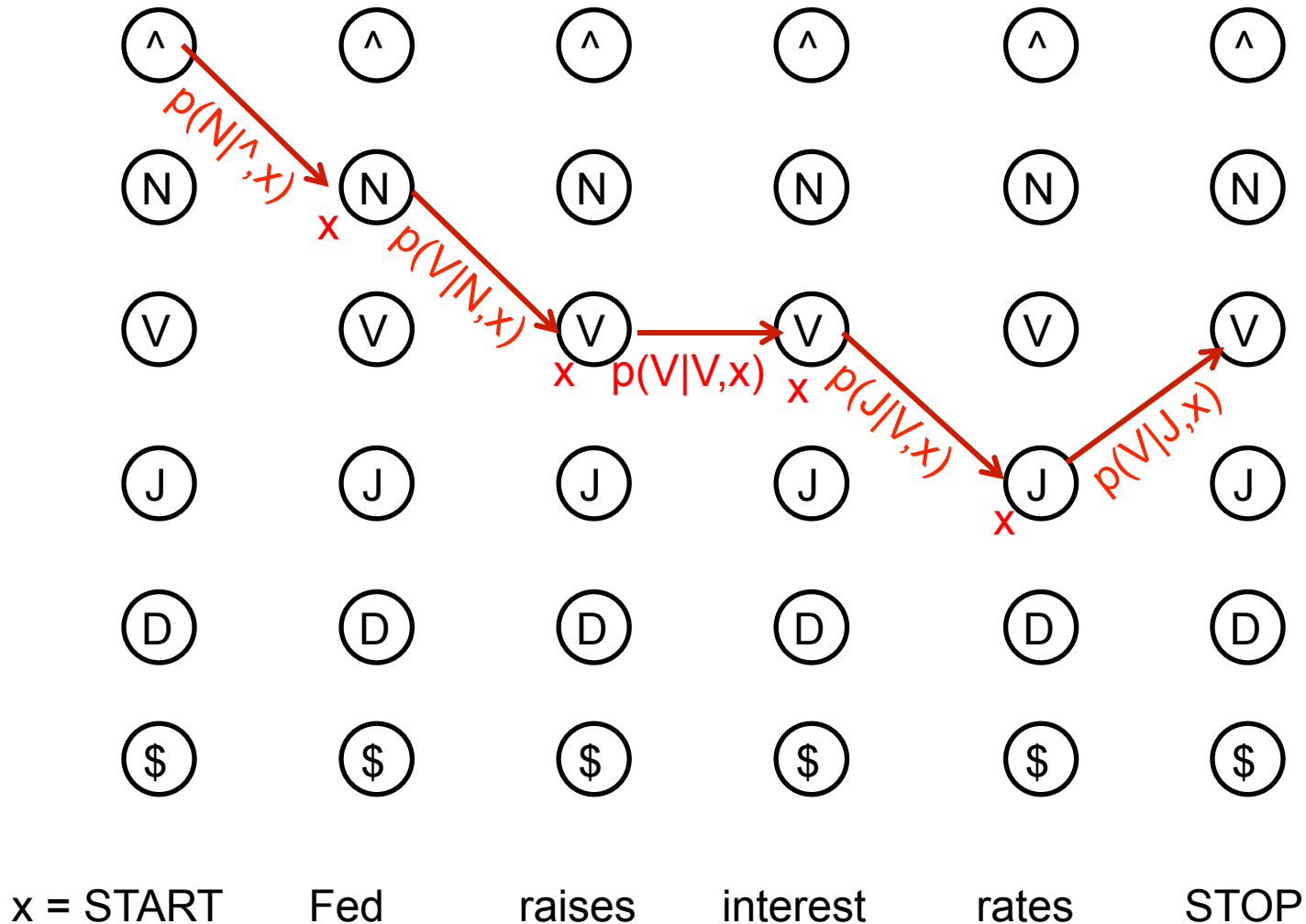
Challenge: How to compute argmax efficiently?

Decoding

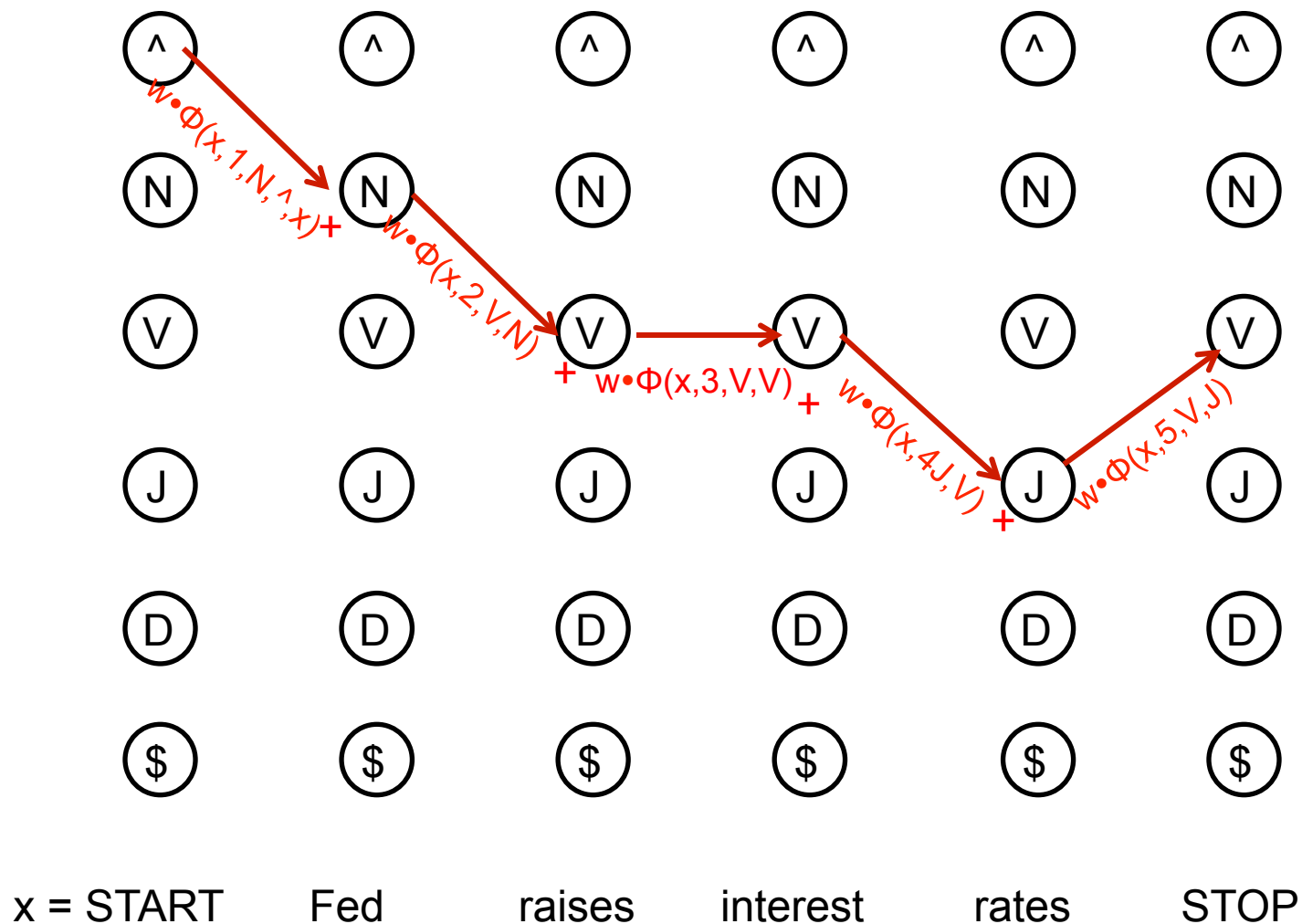
- **Linear Perceptron** $s^* = \arg \max_s w \cdot \Phi(x, s) \cdot \theta$
 - Features must be local, for $x=x_1 \dots x_m$, and $s=s_1 \dots s_m$

$$\Phi(x, s) = \sum_{j=1}^m \phi(x, j, s_{j-1}, s_j)$$

The MEMM State Lattice / Trellis (repeat)



The Perceptron State Lattice / Trellis



Decoding

- **Linear Perceptron** $s^* = \arg \max_s w \cdot \Phi(x, s) \cdot \theta$

- Features must be local, for $x=x_1 \dots x_m$, and $s=s_1 \dots s_m$

$$\Phi(x, s) = \sum_{j=1}^m \phi(x, j, s_{j-1}, s_j)$$

- Define $\pi(i, s_i)$ to be the max score of a sequence of length i ending in tag s_i

$$\pi(i, s_i) = \max_{s_{i-1}} w \cdot \phi(x, i, s_{i-1}, s_i) + \pi(i-1, s_{i-1})$$

- **Viterbi algorithm (HMMs):**

$$\pi(i, s_i) = \max_{s_{i-1}} e(x_i | s_i) q(s_i | s_{i-1}) \pi(i-1, s_{i-1})$$

- **Viterbi algorithm (Maxent):**

$$\pi(i, s_i) = \max_{s_{i-1}} p(s_i | s_{i-1}, x_1 \dots x_m) \pi(i-1, s_{i-1})$$

Overview: Accuracies

- Roadmap of (known / unknown) accuracies:
 - Most freq tag: ~90% / ~50%
 - Trigram HMM: ~95% / ~55%
 - TnT (HMM++): 96.2% / 86.0%
 - Maxent $P(s_i|x)$: 96.8% / 86.8%
 - MEMM tagger: 96.9% / 86.9%
 - Perceptron 96.7% / ??

- Upper bound: ~98%

Probabilistic Models

(Unstructured) categorization:

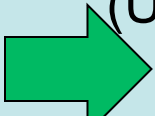
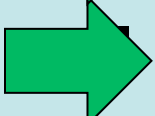
- Naïve Bayes

Structured prediction:

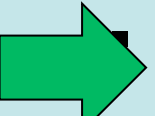
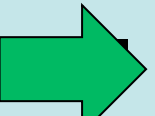
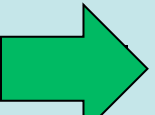
- HMMs
- PCFG Models
- IBM Models

Feature-rich / (Log)-linear Models

(Unstructured) categorization:

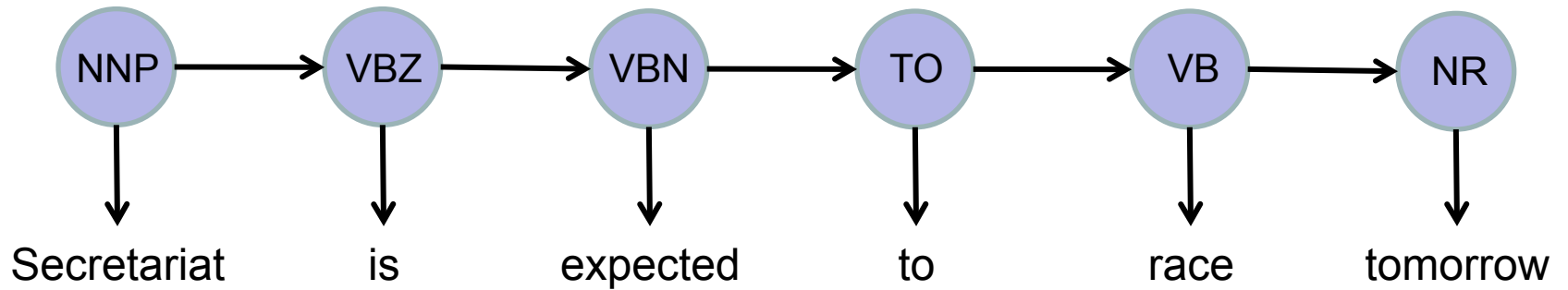
-  Perceptron
-  Maximum Entropy

Structured prediction:

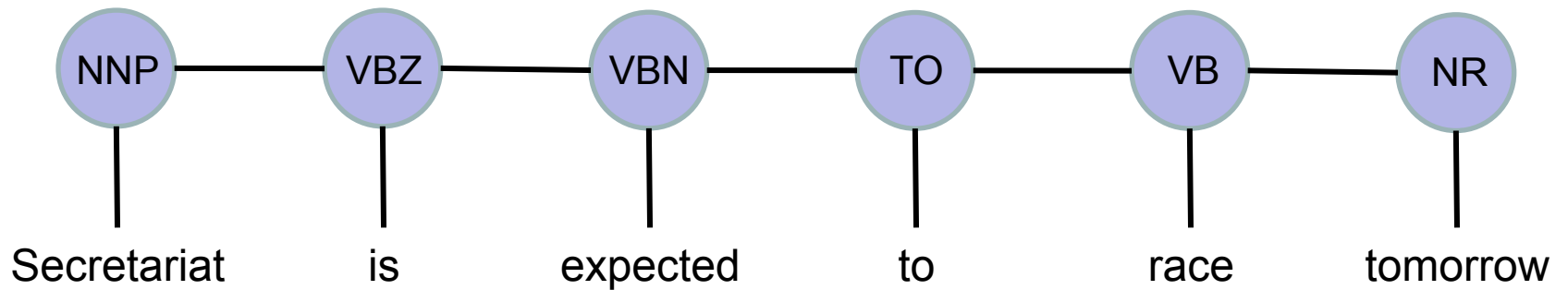
-  Perceptron for Structured Prediction
-  MEMM (Maximum Entropy Markov Model)
-  CRF (Conditional Random Fields)

MEMM v.s. CRF (Conditional Random Fields)

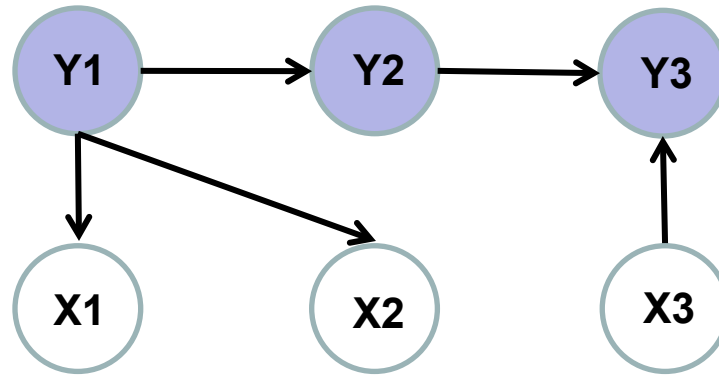
MEMM



CRF

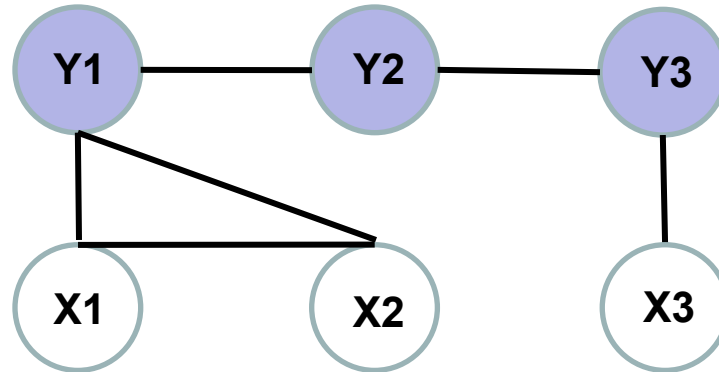


Graphical Models



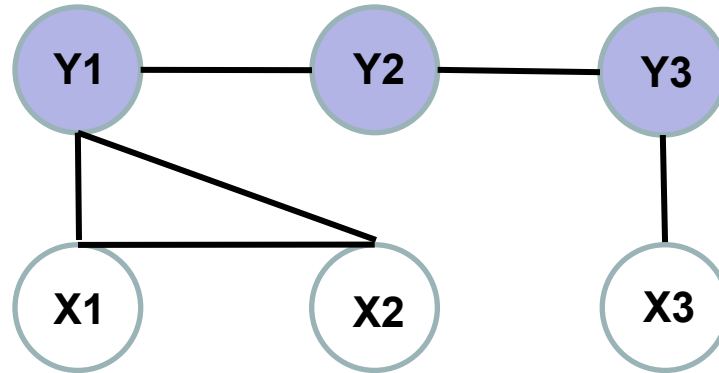
- Conditional probability for each node
 - e.g. $p(Y3 | Y2, X3)$ for Y3
 - e.g. $p(X3)$ for X3
- Conditional independence
 - e.g. $p(Y3 | Y2, X3) = p(Y3 | Y1, Y2, X1, X2, X3)$
- Joint probability of the entire graph
 - = product of conditional probability of each node

Undirected Graphical Model Basics



- Conditional independence
 - e.g. $p(Y3 \mid \text{all other nodes}) = p(Y3 \mid Y3\text{' neighbor})$
- No conditional probability for each node
- Instead, “**potential function**” for each **clique**
 - e.g. $\phi(X1, X2, Y1)$ or $\phi(Y1, Y2)$
- Typically, log-linear potential functions
 - ➔ $\phi(Y1, Y2) = \exp \sum_k w_k f_k(Y1, Y2)$

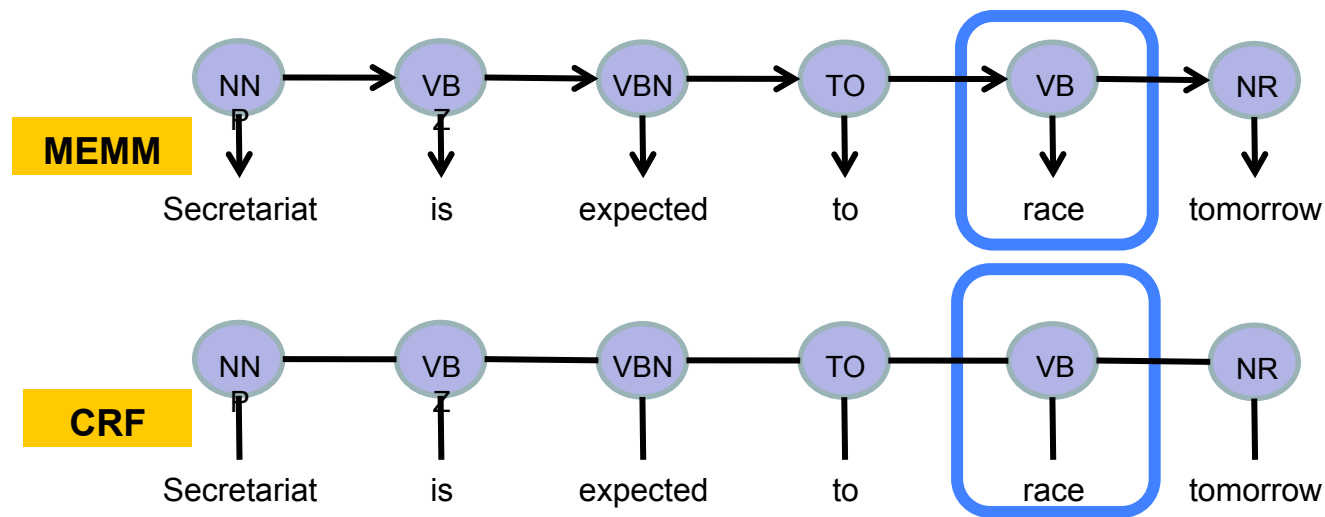
Undirected Graphical Model Basics



- Joint probability of the entire graph

$$P(\vec{Y}) = \frac{1}{Z} \prod_{\text{clique } C} \varphi(\vec{Y}_C)$$

$$Z = \sum_{\vec{Y}} \prod_{\text{clique } C} \varphi(\vec{Y}_C)$$



MEMM

Directed graphical model

CRF

Undirected graphical model

“Discriminative” or “Conditional” models
 → conditional probability $p(\text{tags} \mid \text{words})$

Probability is defined for each slice =

$$P(\text{tag}_i \mid \text{tag}_{i-1}, \text{word}_i)$$

or

$$p(\text{tag}_i \mid \text{tag}_{i-1}, \text{all words})$$

Instead of probability, **potential (energy function)** is defined for each slide =

$$\phi(\text{tag}_i, \text{tag}_{i-1}) * \phi(\text{tag}_i, \text{word}_i)$$

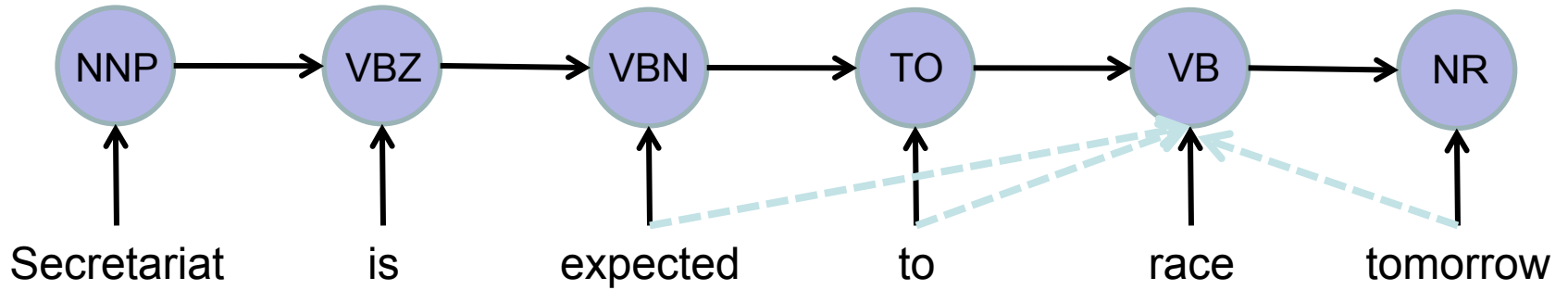
or

$$\phi(\text{tag}_i, \text{tag}_{i-1}, \text{all words}) * \phi(\text{tag}_i, \text{all words})$$

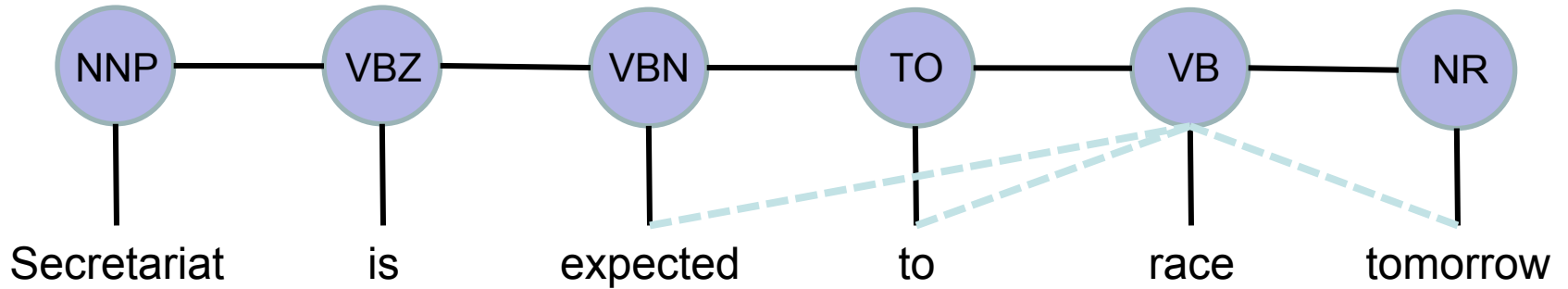
→ Can incorporate long distance features

MEMM v.s. CRF

MEMM



CRF



Conditional Random Fields (CRFs)

[Lafferty, McCallum, Pereira 01]

- Maximum entropy (logistic regression)

Sentence: $x = x_1 \dots x_m$

$$p(y|x; w) = \frac{\exp(w \cdot \phi(x, y))}{\sum_{y'} \exp(w \cdot \phi(x, y'))}$$

Tag Sequence: $y = s_1 \dots s_m$

- Learning:** maximize the (log) conditional likelihood of training data $\{(x_i, y_i)\}_{i=1}^n$

$$\frac{\partial}{\partial w_j} L(w) = \sum_{i=1}^n \left(\phi_j(x_i, y_i) - \sum_y p(y|x_i; w) \phi_j(x_i, y) \right) - \lambda w_j$$

- Computational Challenges?**

- Most likely tag sequence, normalization constant, gradient

Decoding

$$s^* = \arg \max_s p(s|x; w)$$

■ CRFs

- Features must be local, for $x=x_1 \dots x_m$, and $s=s_1 \dots s_m$

$$p(s|x; w) = \frac{\exp(w \cdot \Phi(x, s))}{\sum_{s'} \exp(w \cdot \Phi(x, s'))} \quad \Phi(x, s) = \sum_{j=1}^m \phi(x, j, s_{j-1}, s_j)$$

$$\arg \max_s \frac{\exp(w \cdot \Phi(x, s))}{\sum_{s'} \exp(w \cdot \Phi(x, s'))} = \arg \max_s \exp(w \cdot \Phi(x, s))$$

$$= \arg \max_s w \cdot \Phi(x, s)$$

■ Same as Linear Perceptron!!!

$$\pi(i, s_i) = \max_{s_{i-1}} \phi(x, i, s_{i-1}, s_i) + \pi(i-1, s_{i-1})$$

CRFs: Computing Normalization*

$$p(s|x; w) = \frac{\exp(w \cdot \Phi(x, s))}{\sum_{s'} \exp(w \cdot \Phi(x, s'))} \quad \Phi(x, s) = \sum_{j=1}^m \phi(x, j, s_{j-1}, s_j)$$

$$\begin{aligned} \sum_{s'} \exp(w \cdot \Phi(x, s')) &= \sum_{s'} \exp\left(\sum_j w \cdot \phi(x, j, s_{j-1}, s_j)\right) \\ &= \sum_{s'} \prod_j \exp(w \cdot \phi(x, j, s_{j-1}, s_j)) \end{aligned}$$

Define $norm(i, s_i)$ to sum of scores for sequences ending in position i

$$norm(i, y_i) = \sum_{s_{i-1}} \exp(w \cdot \phi(x, i, s_{i-1}, s_i)) norm(i-1, s_{i-1})$$

- **Forward Algorithm! Remember HMM case:**

$$\alpha(i, y_i) = \sum_{y_{i-1}} e(x_i | y_i) q(y_i | y_{i-1}) \alpha(i-1, y_{i-1})$$

- Could also use backward?

CRFs: Computing Gradient*

$$p(s|x; w) = \frac{\exp(w \cdot \Phi(x, s))}{\sum_{s'} \exp(w \cdot \Phi(x, s'))} \quad \Phi(x, s) = \sum_{j=1}^m \phi(x, j, s_{j-1}, s_j)$$

$$\frac{\partial}{\partial w_j} L(w) = \sum_{i=1}^n \left(\Phi_j(x_i, s_i) - \sum_s p(s|x_i; w) \Phi_j(x_i, s) \right) - \lambda w_j$$

$$\begin{aligned} \sum_s p(s|x_i; w) \Phi_j(x_i, s) &= \sum_s p(s|x_i; w) \sum_{j=1}^m \phi_k(x_i, j, s_{j-1}, s_j) \\ &= \sum_{j=1}^m \sum_{a,b} \sum_{s: s_{j-1}=a, s_b=b} p(s|x_i; w) \phi_k(x_i, j, s_{j-1}, s_j) \end{aligned}$$

- Need forward and backward messages

See notes for full details!

Overview: Accuracies

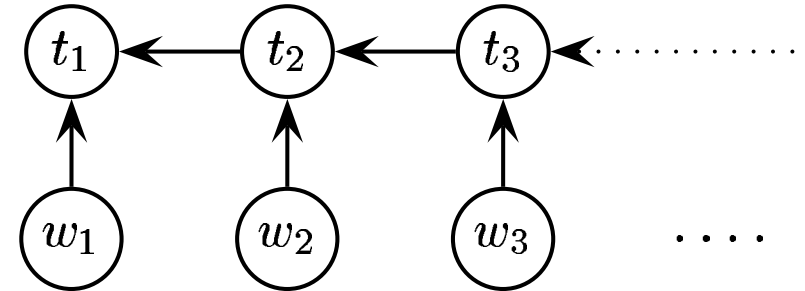
- Roadmap of (known / unknown) accuracies:
 - Most freq tag: ~90% / ~50%
 - Trigram HMM: ~95% / ~55%
 - TnT (HMM++): 96.2% / 86.0%
 - Maxent $P(s_i|x)$: 96.8% / 86.8%
 - MEMM tagger: 96.9% / 86.9%
 - Perceptron 96.7% / ??
 - CRF (untuned) 95.7% / 76.2%
- Upper bound: ~98%

Cyclic Network [Toutanova et al 03]

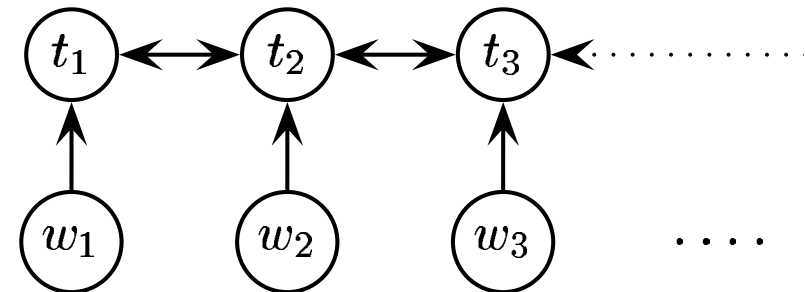
- Train two MEMMs, multiple together to score
- And be very careful
 - Tune regularization
 - Try lots of different features
 - See paper for full details



(a) Left-to-Right CMM



(b) Right-to-Left CMM



(c) Bidirectional Dependency Network

Figure 1: Dependency networks: (a) the (star first-order CMM), (b) the (reversed) right-to-left the bidirectional dependency network.

Overview: Accuracies

- Roadmap of (known / unknown) accuracies:
 - Most freq tag: ~90% / ~50%
 - Trigram HMM: ~95% / ~55%
 - TnT (HMM++): 96.2% / 86.0%
 - Maxent $P(s_i|x)$: 96.8% / 86.8%
 - MEMM tagger: 96.9% / 86.9%
 - Perceptron 96.7% / ??
 - CRF (untuned) 95.7% / 76.2%
 - Cyclic tagger: 97.2% / 89.0%
 - Upper bound: ~98%