

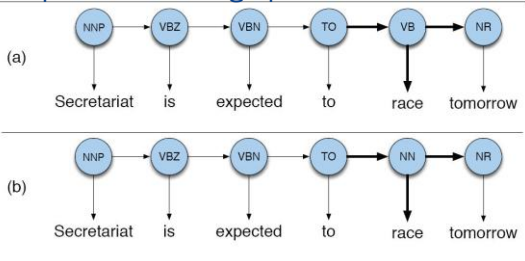
Classification

- y : random variable for prediction (output)
- x : random variable for observation (input)
- Training Data = Collection of (x, y) pairs
- Machine Learning = Given the training data, learn a mapping function $f(x) = y$ that can map input variables to output variables
- Binary classification
- Multiclass classification

Sequence Tagging

- y : **A sequence of** random variables for prediction (output)
- x : **A sequence of** random variables for observation (input)
- Training Data = Collection of (x, y) pairs
- Machine Learning = Given the training data, learn a mapping function $f(x) = y$ that can map input variables to output variables
- Binary classification
- Multiclass classification

Hidden Markov Model (HMM) represented as a graphical model

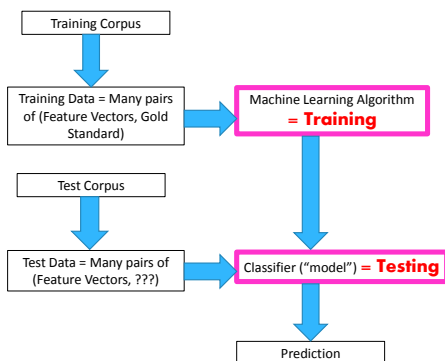
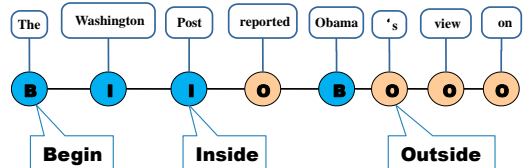


- Note that in this representation, the number of nodes (states) = the length of the word sequence.

Sequence Tagging

- Prediction using **BIO** tagging

<The Washington Post> reported <Obama>'s view on the oil crisis.



- Classifier
 - Maximum Entropy (MaxEnt)
 - Naïve Bayes
- Sequence Tagger
 - Hidden Markov Models (HMMs)
 - Maximum Entropy Markov Models (MEMMs)
 - Conditional Random Fields (CRFs)

Maximum Entropy (MaxEnt) Models

Maximum Entropy (MaxEnt) Models

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

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

- The feature vector representation may include redundant and overlapping features

(slide modified from Dan Klein's)

Training Maximum Entropy (MaxEnt)

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

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

← Make positive
← Normalize

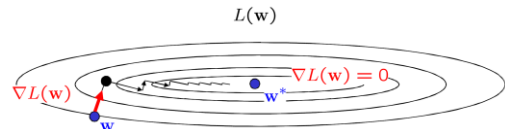
- Maximize the (log) conditional likelihood of training data

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

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

(slide modified from Dan Klein's)

Convex Optimization for Training



- 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:
 - evaluate the function at current ‘w’
 - evaluate its derivative at current ‘w’

(slide modified from Dan Klein's)

Training Maximum Entropy

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

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

Total count of feature n in correct candidates

Expected count of feature n in predicted candidates

(slide modified from Dan Klein's)

Training with Regularization

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

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

Big weights are bad

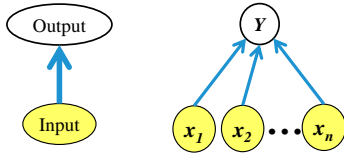
Total count of feature n in correct candidates

Expected count of feature n in predicted candidates

(slide modified from Dan Klein's)

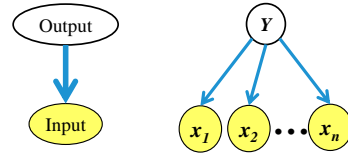
Graphical Representation of MaxEnt

$$P(y|x, w) = \frac{\exp(w^T f(y))}{\sum_{y'} \exp(w^T f(y'))}$$

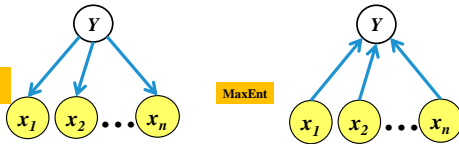


Graphical Representation of Naïve Bayes

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



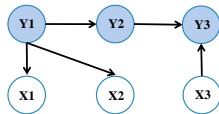
• Consult AI text book for more details



Naïve Bayes Classifier	Maximum Entropy Classifier
"Generative" models → $p(\text{input} \text{output})$ → For instance, for text categorization, $P(\text{words} \text{category})$ → Waste energy on generating input (which we don't need to generate during test) → 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	"Discriminative" models → $p(\text{output} \text{input})$ → For instance, for text categorization, $P(\text{category} \text{words})$ → Focusing only 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 → Can handle redundant and overlapping features

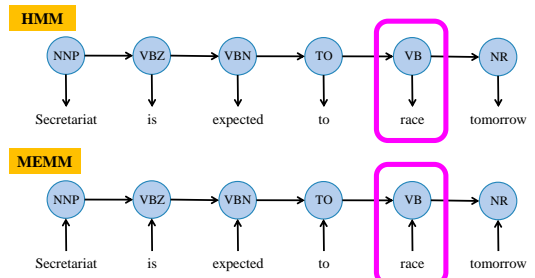
Sequence Tagging with HMM / MEMM / CRF

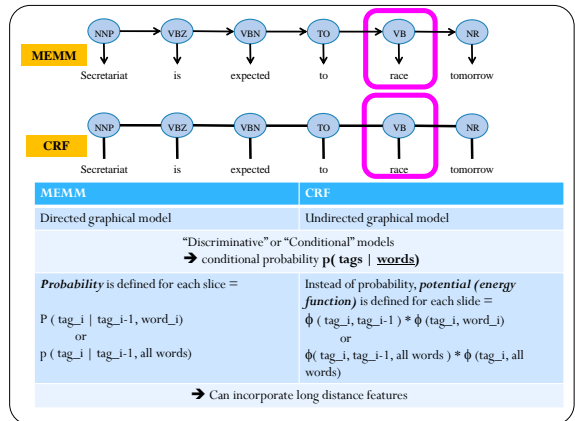
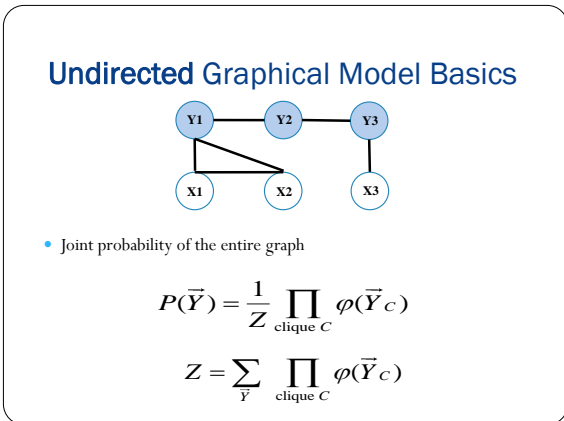
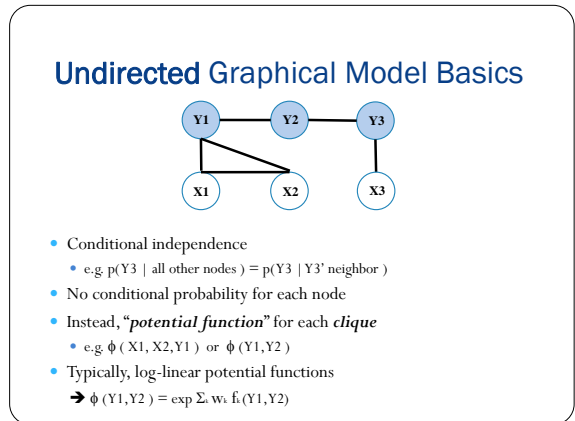
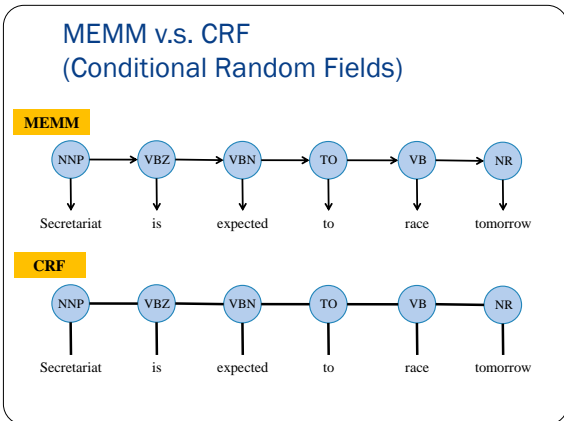
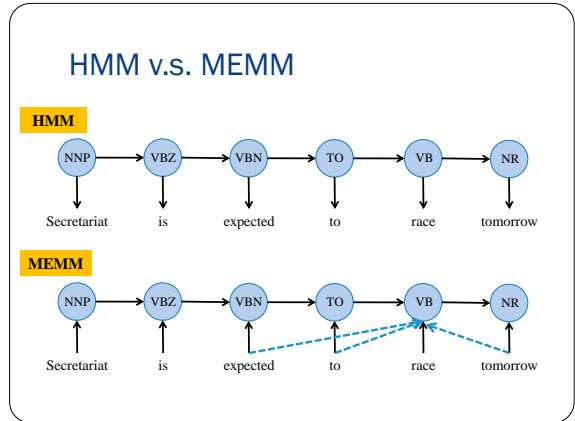
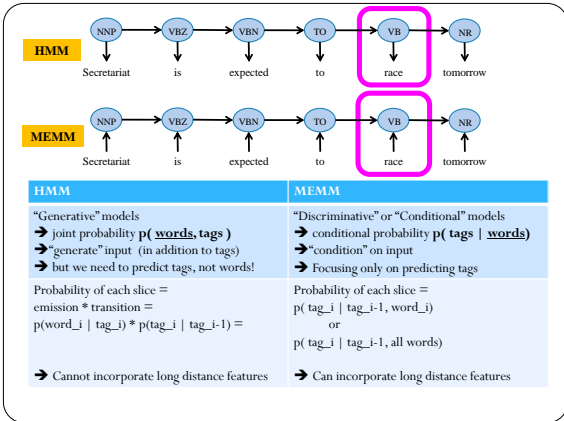
Graphical Model Basics



- Conditional probability for each node
 - e.g. $p(Y_3 | Y_2, X_3)$ for Y_3
 - e.g. $p(X_3)$ for X_3
- Conditional independence
 - e.g. $p(Y_3 | Y_2, X_3) = p(Y_3 | Y_1, Y_2, X_1, X_2, X_3)$
- Joint probability of the entire graph
= product of conditional probability of each node

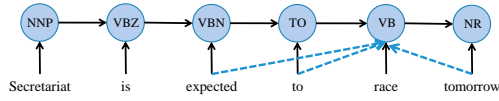
HMM v.s. MEMM (Maximum Entropy Markov Models)



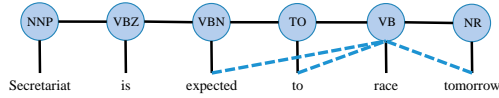


MEMM v.s. CRF

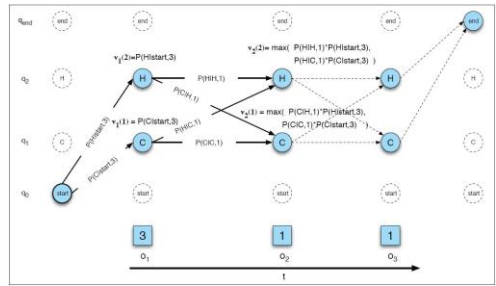
MEMM



CRF



Inference (Viterbi)



Objective function for training

Given the training data $D = \{x^{(j)}, y^{(j)}\}_{j=1}^N$
 and $p(y | x) = \frac{1}{Z(x)} \exp \lambda \cdot \mathbf{F}(y, x)$

Objective function :
 conditional likelihood $L(\lambda) = L(\lambda | D) = P(D | \lambda) = \prod_j p(y^{(j)} | x^{(j)})$
 equiv. to optimize $l(\lambda) = \log L(\lambda) = \sum_j \log p(y^{(j)} | x^{(j)})$

$$\begin{aligned}
 l(\lambda) &= \sum_j \log p(y^{(j)} | x^{(j)}) = \sum_j \log \frac{1}{Z(x^{(j)})} \exp \lambda \cdot \mathbf{F}(y^{(j)}, x^{(j)}) \\
 &= \sum_j \lambda \cdot \mathbf{F}(y^{(j)}, x^{(j)}) - \log Z(x^{(j)}) \\
 &= \sum_j \left(\lambda \cdot \mathbf{F}(y^{(j)}, x^{(j)}) - \log \sum_y \exp \lambda \cdot \mathbf{F}(y, x^{(j)}) \right)
 \end{aligned}$$

nasty!

CRFs Software:

- Mallet (<http://mallet.cs.umass.edu/>),
- CRF++ (<http://crfpp.sourceforge.net/>),
- CRF (<http://crf.sourceforge.net/>)