

Conditional Random Fields and structured prediction

CS-585

Natural Language Processing

Sonjia Waxmonsky

HMMs, MEMMs and CRFs

- Hidden Markov Models (HMM)
 - Generative sequence labeling models
- Maximum Entropy Markov Models (MEMM)
 - Like HMMs, but discriminative
 - Logistic regression
- Conditional Random Fields (CRF)
 - Like MEMMs, but in structured prediction paradigm

MAXIMUM ENTROPY MARKOV MODELS

Maximum Entropy

- In NLP, logistic regression models are often called *maximum entropy* (or *maxent*) models
- Idea
 - When selecting a probability distribution to model observed data, we want the distribution to model the statistics of the data
 - Model should be as uncertain as possible, while still capturing the data
 - Uncertainty=entropy, so select the model with maximal entropy subject to data-fitting constraints

Maximum Entropy

- In a supervised learning context, the data fitting constraints have to do with the frequencies of feature functions $f_k(X, Y)$ in the training data:

$$\sum_n f_k(X_n, Y_n) P(X_n, Y_n) = c_k$$

- It can be shown that the distribution $P(Y_n|X_n)$ with maximum entropy subject to these constraints has the form

$$P(Y_n|X_n) = \frac{e^{\sum_k \lambda_k f_k(X_n, Y_n)}}{Z}$$

- Where Z is a normalizing constant
- Just logistic regression....

Maximum Entropy Markov Models (MEMMs)

- Cousins of HMMs

- Still based on the Markov assumption:

$$P(T|W) = \prod_i P(T_i | T_{i-k..i-1}, W)$$

- Because they are discriminative, MEMMs cannot be trained in an unsupervised way like HMMs
 - HMMs are generative (model $P(W, T)$)
 - MEMMs are discriminative (model $P(T|W)$)
- Based on maximum entropy (logistic regression) models to predict tag for each word given local context
 - Can be trained using gradient descent or model-specific algorithms (GIS, IIS)

Maximum Entropy Markov Models (MEMMs)

Estimate conditional probability of tags given text

$$P_m(T|W) = \prod_i P(T_i|T_{i-k..i-1}, W)$$

(Markov assumption)

$$= \prod_i \frac{e^{\sum_k \lambda_k f_k(W, T_{i-k..i})}}{\sum_{t \in \mathcal{T}} e^{\sum_k \lambda_k f_k(W, T_{i-k..i-1}, t)}}$$

Locally normalized

(maxent)

- Can be trained using gradient descent
 - Either maximum-likelihood estimation or penalized likelihood (regularization)
- Similar dynamic programming tricks to HMMs for efficiently computing sum across all labelings

MEMMs: POS feature functions


- Maxent feature functions

Condition	Features
w_i is not rare	$w_i = X$ & $t_i = T$
w_i is rare	X is prefix of w_i , $ X \leq 4$ & $t_i = T$
	X is suffix of w_i , $ X \leq 4$ & $t_i = T$
	w_i contains number & $t_i = T$
	w_i contains uppercase character & $t_i = T$
	w_i contains hyphen & $t_i = T$
$\forall w_i$	$t_{i-1} = X$ & $t_i = T$
	$t_{i-2}t_{i-1} = XY$ & $t_i = T$
	$w_{i-1} = X$ & $t_i = T$
	$w_{i-2} = X$ & $t_i = T$
	$w_{i+1} = X$ & $t_i = T$
	$w_{i+2} = X$ & $t_i = T$

Table 1: Features on the current history h_i

MEMMs: POS feature functions

Word	the	stories	about	well-heeled	communities	and	developers
Tag	DT	NNS	IN	JJ	NNS	CC	NNS
Position	1	2	3	4	5	6	7



$w_i = \text{about}$	$\wedge t_i = \text{IN}$
$w_{i-1} = \text{stories}$	$\wedge t_i = \text{IN}$
$w_{i-2} = \text{the}$	$\wedge t_i = \text{IN}$
$w_{i+1} = \text{well-heeled}$	$\wedge t_i = \text{IN}$
$w_{i+2} = \text{communities}$	$\wedge t_i = \text{IN}$
$t_{i-1} = \text{NNS}$	$\wedge t_i = \text{IN}$
$t_{i-2}t_{i-1} = \text{DT NNS}$	$\wedge t_i = \text{IN}$

MEMMs: POS feature functions

Word	the	stories	about	well-heeled	communities	and	developers
Tag	DT	NNS	IN	JJ	NNS	CC	NNS
Position	1	2	3	4	5	6	7

$prefix(w_i) = w$	$\wedge t_i = JJ$
$prefix(w_i) = we$	$\wedge t_i = JJ$
$prefix(w_i) = wel$	$\wedge t_i = JJ$
$prefix(w_i) = well$	$\wedge t_i = JJ$
$suffix(w_i) = d$	$\wedge t_i = JJ$
$suffix(w_i) = ed$	$\wedge t_i = JJ$
$suffix(w_i) = led$	$\wedge t_i = JJ$
$suffix(w_i) = eled$	$\wedge t_i = JJ$
w_i contains hyphen	$\wedge t_i = JJ$

$w_{i-1} = \text{about}$	$\wedge t_i = JJ$
$w_{i-2} = \text{stories}$	$\wedge t_i = JJ$
$w_{i+1} = \text{communities}$	$\wedge t_i = JJ$
$w_{i+2} = \text{and}$	$\wedge t_i = JJ$
$t_{i-1} = \text{IN}$	$\wedge t_i = JJ$
$t_{i-2}t_{i-1} = \text{NNS IN}$	$\wedge t_i = JJ$

Viterbi Tagging for MEMMs

Most probable tag sequence given text:

$$T^* = \operatorname{argmax}_T P_m(T|W)$$

$$= \operatorname{argmax}_T \prod_i P(T_i | T_{i-k..i-1}, W)$$

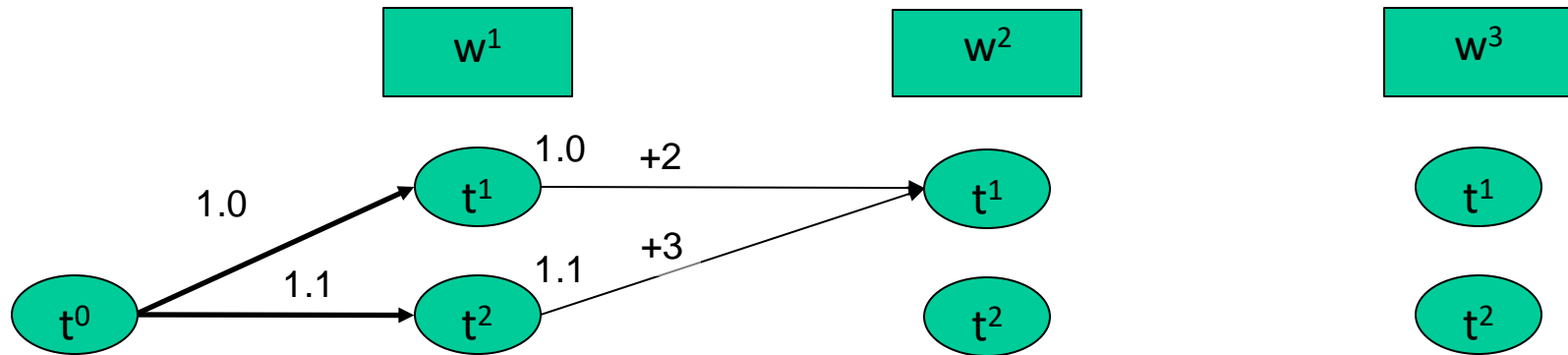
(Markov assumption)

$$= \operatorname{argmax}_T \prod_i \frac{e^{\sum_k \lambda_k f_k(W, T_{i-k..i})}}{\sum_{t \in \mathcal{T}} e^{\sum_k \lambda_k f_k(W, T_{i-k..i-1}, t)}}$$

(maxent)

$$= \operatorname{argmax}_T \sum_i \left(\sum_k \lambda_k f_k(W, T_{i-k..i}) - \log \sum_{t \in \mathcal{T}} e^{\sum_k \lambda_k f_k(W, T_{i-k..i-1}, t)} \right)$$

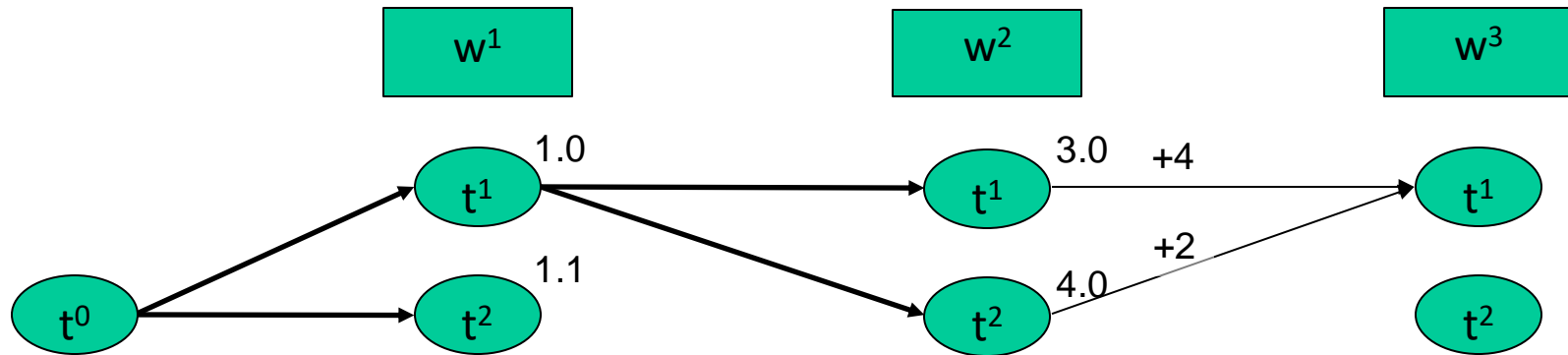
Viterbi algorithm



$$-\log P(t_i | t_{i-1}, w_i, w_{i-1})$$

t _{i-1}	t ⁰		t ¹									t ²								
w _i	w ¹	w ²	w ¹			w ²			w ³			w ¹			w ²			w ³		
w _{i-1}			w ¹	w ²	w ³	w ¹	w ²	w ³	w ¹	w ²	w ³	w ¹	w ²	w ³	w ¹	w ²	w ³	w ¹	w ²	w ³
t _i =t ¹	1.0	1.1	2	4	2	2	3	2	4	5	2	2	2	3	2	4	5	2	4	
t _i =t ²	1.1	2.1	3	5	3	3	2	4	2	4	3	3	3	2	4	2	4	3	4	5

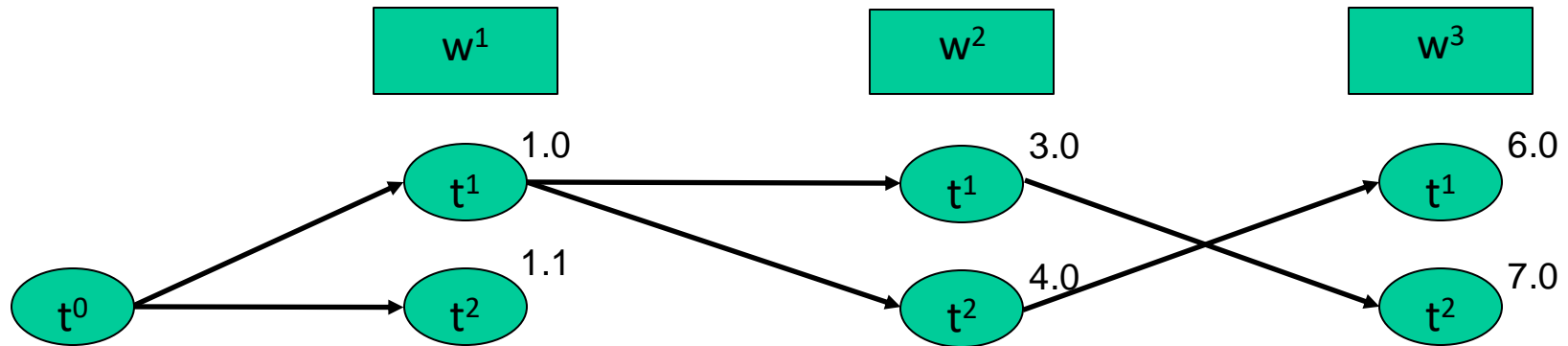
Viterbi algorithm



$$-\log P(t_i | t_{i-1}, w_i, w_{i-1})$$

t_{i-1}	t^0		t^1									t^2								
w_i	w^1	w^2	w^1			w^2			w^3			w^1			w^2			w^3		
w_{i-1}			w^1	w^2	w^3	w^1	w^2	w^3	w^1	w^2	w^3	w^1	w^2	w^3	w^1	w^2	w^3	w^1	w^2	w^3
$t_i = t^1$	1.0	1.1	2	4	2	2	2	3	2	4	5	2	2	2	3	2	4	5	2	4
$t_i = t^2$	1.1	2.1	3	5	3	3	2	4	2	4	3	3	3	2	4	2	4	3	4	5

Viterbi algorithm



$$-\log P(t_i | t_{i-1}, w_i, w_{i-1})$$

t_{i-1}	t^0		t^1									t^2								
w_i	w^1	w^2	w^1			w^2			w^3			w^1			w^2			w^3		
w_{i-1}			w^1	w^2	w^3	w^1	w^2	w^3	w^1	w^2	w^3	w^1	w^2	w^3	w^1	w^2	w^3	w^1	w^2	w^3
$t_i = t^1$	1.0	1.1	2	4	2	2	2	3	2	4	5	2	2	2	3	2	4	5	2	4
$t_i = t^2$	1.1	2.1	3	5	3	3	2	4	2	4	3	3	3	2	4	2	4	3	4	5

Viterbi Algorithm

$D(0, \text{START}) = 0$

for each tag $t \neq \text{START}$ **do:**

$D(0, t) = -\infty$

for $i \leftarrow 1$ **to** N **do:**

for each tag t^j **do:**

$$D(i, t^j) \leftarrow \max_k (D(i-1, t^k) + \log P(t_i | t_{i-1}=t^k, W))$$

$\log P(T|W) = \max_j D(N, t^j)$

MEMM Limitations

- Locally normalized
 - MEMM assumes that the probability distribution over a tagging T can be ***factored*** into the product of conditional probabilities with limited history for each tag location in a sentence
 - This limits the flexibility of the model. Only paths from the same prior state can “compete” against one another for probability mass

CONDITIONAL RANDOM FIELDS

Structured prediction

- CRFs fall into a predictive modeling framework called **structured prediction**
 - When we have some complex structured object over which we want to make predictions...
 - pixels within an image, tag sequences or syntactic trees for a text
 - We try to estimate a *probability distribution* over the *whole* output space, rather than distributions over *subparts*

Structured prediction

Structured prediction for sequence labeling: predict optimal tagging T^* using a scoring function over the output space:

$$\hat{T} = \operatorname{argmax}_T \Psi(T, W)$$

In a probabilistic framework:

$$\hat{T} = \operatorname{argmax}_T P(T|W)$$

Specifically using logistic regression:

$$\hat{T} = \operatorname{argmax}_T \frac{\prod_i e^{\lambda_i f_i(T, W)}}{\sum_{T'} \prod_i e^{\lambda_i f_i(T', W)}}$$

Feature functions for CRFs

- Essentially the same as for MEMMs: can incorporate left and right context for observations (words)
- For tags, the probabilistic framework theoretically could encompass features built on arbitrary tag dependencies
 - E.g., first and last tag in sentence
 - But we need dynamic programming to make inference tractable. Therefore, in practice, tag dependencies are limited to adjacent tags
- “Linear chain CRF”

The partition function

- Remember our expression for the probability of a tag sequence under a CRF:

$$P(T|W) = \frac{\prod_i e^{\lambda_i f_i(T,W)}}{\sum_{T'} \prod_i e^{\lambda_i f_i(T',W)}}$$

The denominator is called the partition function. It involves a sum over all possible tag sequences for the sentence and is expensive to compute

But note that we don't need to compute it in order to predict the best tagging; the numerator is sufficient

Inference in CRFs: the Viterbi algorithm

$$\begin{aligned}\operatorname{argmax}_T P(T|W) &= \operatorname{argmax}_T \frac{\prod_i e^{\lambda_i f_i(T,W)}}{\sum_{T'} \prod_i e^{\lambda_i f_i(T',W)}} \\ &= \operatorname{argmax}_T \prod_i e^{\lambda_i f_i(T,W)} \\ &= \operatorname{argmax}_T \sum_i \lambda_i f_i(T,W)\end{aligned}$$

Notebook – CONLL with CRF-Suite library

Training of CRFs: the forward recurrence

The conditional log likelihood of the training data under the model is

$$\begin{aligned}\log P_m(T^*|W) &= \sum_{k=1}^N \log \frac{\prod_i e^{\lambda_i f_i(T_k^*, W_k)}}{\sum_{T'} \prod_i e^{\lambda_i f_i(T', W_k)}} \\ &= \sum_{k=1}^N \left(\sum_i \lambda_i f_i(T_k^*, W_k) - \log \sum_{T'} \prod_i e^{\lambda_i f_i(T', W_k)} \right)\end{aligned}$$

We want to compute gradients of this so that we can use gradient descent for optimization

Training of CRFs: the forward recurrence

$$\sum_{k=1}^N \left(\sum_i \lambda_i f_i(T_k^*, W_k) - \log \sum_{T'} \prod_i e^{\lambda_i f_i(T', W_k)} \right)$$

Sum over all
training data

Sum of feature potentials;
easy to compute

Partition function; expensive
to compute

Use dynamic programming!

Training of CRFs: the forward recurrence

Define forward variable $a_i(t^j)$ as the sum of scores of all paths ending up with tag t^j at word w_i :

$$a_i(t^j) = \sum_{t' \in \mathcal{T}} \left(e^{\lambda_k f_k(t', t^j, w)} \times a_{i-1}(t') \right)$$

Sum over all
possible previous
tags

Score for transitioning from
previous tag to current tag

Forward variable for
previous tag with
previous word

MEMMs and CRFs: Key Points

- MEMMs and CRFs incorporate complex *features* for sequence labeling in a probabilistic framework
- MEMMs and CRFs are based on logistic regression (a.k.a. maximum entropy)
- CRFs improve on MEMMs by using globally, rather than locally normalized probabilities
- CRFs excel at incorporating global constraints on well-formedness of label sequences
- As with logistic regression, we can apply **L1 and L2 regularization**