

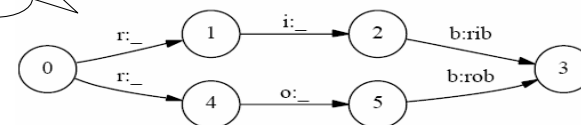
Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data

John Lafferty, Andrew McCallum, Fernando Pereira

Presenter: Yejin Choi

Label Bias Problem

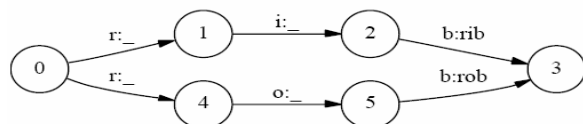
arcs: observations (X)
nodes: outputs (Y)



S1	Herbivores (-)	like (r)	merino (i)	sheep (b)
	NNS (0)	IN (1)	JJ (2)	NN (3)
S2	Carnivores (-)	like (r)	eating (o)	sheep (b)
	NNS (0)	VB (4)	VBG (5)	NN (3)

- Suppose [NNS => VB] transition more frequent than [NNS => IN]
- Suppose from [VB], only [VB => VBG] transition is possible
- ➔ now, What is $P(Y_i = \text{VBG} \mid Y_{i-1} = \text{VB}, X_i = \text{merino})$???
- Recall MEMM models $P(Y_i \mid Y_{i-1}, X_i)$

Label Bias Problem



Herbivores (h)	like (r)	merino (i)	sheep (b)
NNS (0)	IN (1)	JJ (2)	NN (3)
Carnivores (c)	like (r)	eating (o)	sheep (b)
NNS (0)	VB (4)	VBG (5)	NN (3)

So, how do we fix this nonsense ?

$$P(Y_i = \text{VBG} \mid Y_{i-1} = \text{VB}, X_i = \text{merino}) = 1 \text{ for any } X_i$$

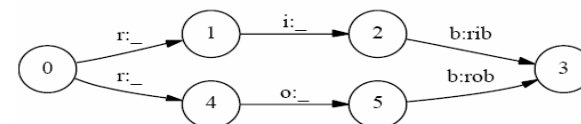
➔ Do not normalize on each node!

Instead, normalize over the entire sequence.

This motivates the "global normalization" scheme of CRFs.

Other approaches: Cohen and Carvalho (2005), Sutton and McCallum (2005)

Label Bias Problem



Herbivores (h)	like (r)	merino (i)	sheep (b)
NNS (0)	IN (1)	JJ (2)	NN (3)
Carnivores (c)	like (r)	eating (o)	sheep (b)
NNS (0)	VB (4)	VBG (5)	NN (3)

Wait, what about HMMs?

- Recall HMM models $P(X_i \mid Y_i)$ and $P(Y_i \mid Y_{i-1})$

$$(\text{so that } P(X_i \mid Y_i, Y_{i-1}) \cdot P(Y_i \mid Y_{i-1}) = P(X_i, Y_i \mid Y_{i-1}))$$

then we can get $P(\text{merino} \mid \text{VBG}) = 0$!

Label Bias v.s. Observation Bias

Label Bias	(Bottou, 1991; Lafferty et. al. 2001)	Previous state explains current state so well (hence ignoring observation)
Observation Bias	(Klein and Manning 2002)	Observation explains current state so well (hence ignoring state transition)

e.g.

All the indexes above	HMM	MEMM
PDT DT NNS VBD	-0.0	-1.3
DT DT NNS VBD	-5.4	-0.3

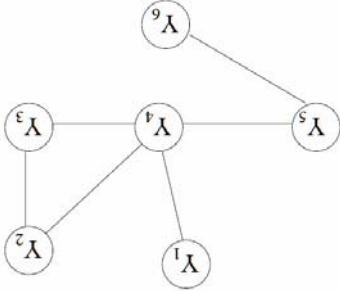
➔ Both are to do with local conditional normalization.

Markov Random Field

Just for now, forget about conditioning on $X_{...}$

Let $G = (Y, E)$ be a graph where each vertex Y_v is a random variable
Suppose $P(Y_v | \text{all other } Y) = P(Y_v | \text{neighbors}(Y_v))$ then Y is a random field

Example:



- $P(Y_5 | \text{all other } Y) = P(Y_5 | Y_4, Y_6)$
- But, how do we compute the global joint distribution $P(Y)$ out of this ?
- Besides, We don't want to compute $P(Y_1 | \text{neighbors}(Y_1))$!!!

Hammersley-Clifford theorem (1971)

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

where $Z = \sum_Y \prod_C \phi(Y_C)$

1. Given MRF $G=(Y,E)$ such that $P(Y_i | Y \setminus Y_i) = P(Y_i | \text{nbr}(Y_i))$
2. Given $\phi(Y_c)$ for \forall clique C in G , such that $\phi(Y_c) >= 0$

- cliques may overlap.
- cliques may not be maximal.

➔ this implies we don't need to compute $P(Y_i | \text{nbr}(Y_i))$ to get $P(Y)$!

So, Let's normalize globally! (But, how?)

- What we need is the global joint distribution. i.e., $p(Y | x)$
 - where $y = (y_1, \dots, y_n)$ and $x = (x_1, \dots, x_n)$
 - we do not want distributions on individual node. i.e., $p(Y_i | x_j)$
 - Instead, we want non-probabilistic *potential* function. i.e., $\phi(Y_i, x_j)$
- $p(y | x) \approx g(\phi(y_1, x_1), \dots, \phi(y_n, x_n))$????????

- Problem with directed graphs (like Bayesian Network)
 - a probability distribution should be given for each node
 - then, the joint probability $p(y) = \prod_i p(y_i | \text{parents}(y_i))$

➔ btw, what is $\text{parents}(y_i)$ for MEMM ?

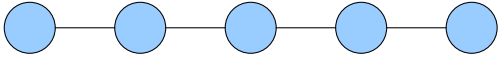
- Markov Random Field ! (= Markov Network, Random Field)

Definition of CRFs

Definition. Let $G = (V, E)$ be a graph such that $\mathbf{Y} = (\mathbf{Y}^v)_{v \in V}$, so that \mathbf{Y} is indexed by the vertices of G . Then (\mathbf{X}, \mathbf{Y}) is a *conditional random field* in case, when conditioned on \mathbf{X} , the random variables \mathbf{Y}^v obey the Markov property with respect to the graph: $p(\mathbf{Y}^v | \mathbf{X}, \mathbf{Y}_w, w \neq v) = p(\mathbf{Y}^v | \mathbf{X}, \mathbf{Y}_w, w \sim v)$, where $w \sim v$ means that w and v are neighbors in G .

- Almost identical to the definition of MRFs, except CRFs condition on \mathbf{X} .
- and notice that \mathbf{X} is not part of V .
- Hammersley-Clifford theorem can be extended to the conditional cases (on \mathbf{X}) as well!

$P(y | x)$ for linear-chain CRFs



Let $\phi(y_i, x) = \exp \sum_k \lambda_k f_k(y_i, x)$
 $\phi(y_i, y_{i+1}, x) = \exp \sum_k \lambda_k f_k(y_i, y_{i+1}, x)$

Then,

$$p(y | x) = \frac{1}{Z(x)} \prod_i \left(\phi(y_i, x) \phi(y_i, y_{i+1}, x) \right) = \frac{1}{Z(x)} \exp \left(\sum_{i,k} \lambda_k f_k(y_i, x) + \sum_{i,k} \lambda_k f_k(y_i, y_{i+1}, x) \right) = \frac{1}{Z(x)} \exp \left(\lambda \bullet \mathbf{F}(y, x) \right)$$

where $Z(x) = \sum_y \exp \left(\lambda \bullet \mathbf{F}(y, x) \right)$

$$p_\theta(y | x) = \frac{1}{Z_\theta(x)} \exp \left(\sum_{i \in E} \lambda_i f_i(e, y|_i, x) + \sum_{i \in V} \lambda_i g_i(v, y|_i, x) \right)$$

Can use any other non negative function, but exp function will give us a nice continuous convex objective function later.

Objective function for training

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

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

$$l(\lambda) = \sum_t \log p(y^{(t)} | x^{(t)}) = \sum_t \log \frac{1}{Z(x^{(t)})} \exp \lambda \bullet \mathbf{F}(y^{(t)}, x^{(t)})$$

$$= \sum_t \lambda \bullet \mathbf{F}(y^{(t)}, x^{(t)}) - \log Z(x^{(t)})$$

$$= \sum_t \left(\lambda \bullet \mathbf{F}(y^{(t)}, x^{(t)}) - \log \sum_y \exp \lambda \bullet \mathbf{F}(y^{(t)}, x^{(t)}) \right)$$

Objective function for training

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

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

$$l(\lambda) = \sum_t \log p(y^{(t)} | x^{(t)}) = \sum_t \log \frac{1}{Z(x^{(t)})} \exp \lambda \bullet \mathbf{F}(y^{(t)}, x^{(t)})$$

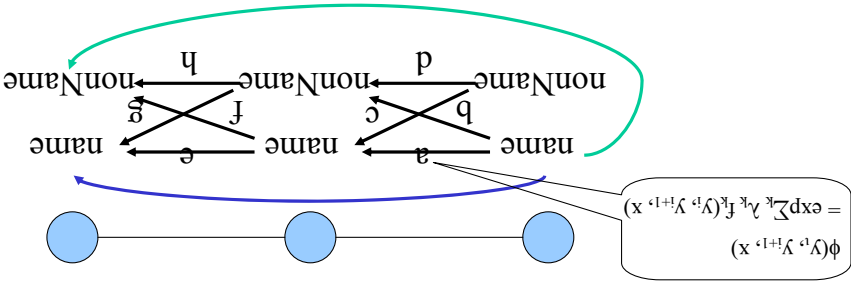
$$= \sum_t \lambda \bullet \mathbf{F}(y^{(t)}, x^{(t)}) - \log Z(x^{(t)})$$

$$= \sum_t \left(\lambda \bullet \mathbf{F}(y^{(t)}, x^{(t)}) - \log \sum_y \exp \lambda \bullet \mathbf{F}(y^{(t)}, x^{(t)}) \right)$$

nasty!

Computing $Z(x)!$... for linear-chain CRFs

$$l(\lambda) = \sum_j \log p(y^{(j)} | x^{(j)}) = \sum_j \left(\lambda \bullet F(y^{(j)}, x^{(j)}) - \log \sum_y \exp \lambda \bullet F(y, x^{(j)}) \right)$$



$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \times \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{bmatrix}$$

Parameter Estimation for CRFs

- Iterative scaling algorithms
 - Generalized iterative scaling (GIS)
 - Improved iterative scaling (IIS)

- Gradient descent methods
 - CG : conjugate gradient
 - L-BFGS : limited memory Newton's method

• All of these maximize the (conditional) likelihood.
 • Possible to train discriminatively with **voted perceptron**.
 (Collins 2002)

Parameter Estimation for CRFs

➔ how to compute $\text{argmax}_{\lambda} l(\lambda)$?

- Iterative scaling algorithms
 - Generalized iterative scaling (GIS)
 - Improved iterative scaling (IIS)
 - easy to implement
 - both really slow to converge
- Gradient descent methods
 - CG : conjugate gradient
 - L-BFGS : limited memory Newton's method
 - much harder to implement, but lots of code available
 - ➔ you only need to provide $l(\lambda)$ and $l'(\lambda)$
 - much faster to converge

Sha and Pereira (2003)

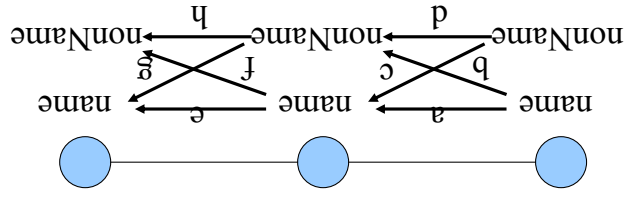
training method	time	F score	\mathcal{L}_{λ}^*
Precond. CG	130	94.19%	-2968
Mixed CG	540	94.20%	-2990
Plain CG	648	94.04%	-2967
L-BFGS	84	94.19%	-2948
GIS	3700	93.55%	-5668

Table 3: Runtime for various training methods in minutes, 375k examples

A universal component for SVM, HMM, MEMM, etc

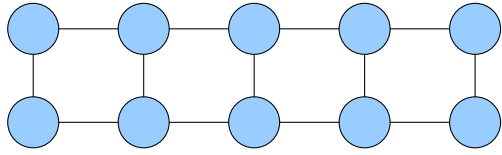
$$\text{argmax}_y P(y | x) \dots \text{ for linear-chain CRFs}$$

$$\begin{aligned} \text{argmax}_y p(y | x) &= \text{argmax}_y \log p(y | x) \\ &= \text{argmax}_y \left(\lambda \bullet F(y, x) - \log Z(x) \right) \\ &= \text{argmax}_y \lambda \bullet F(y, x) \end{aligned}$$

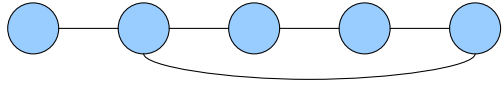


Other CRFs

Factorial CRFs (Sutton et al., 2004)

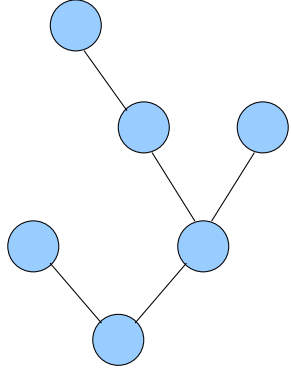


Skip-chain CRFs (Sutton and McCallum, 2004)



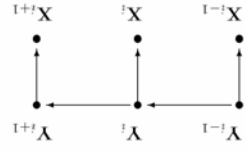
Tree CRFs

(Cohn and Blunsom, 2005)

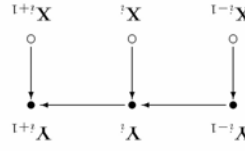


Graphical representation

HMM



MEMM



CRF

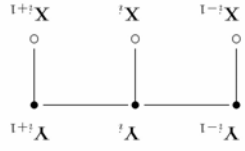
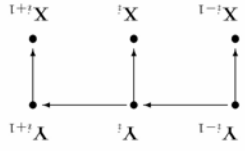
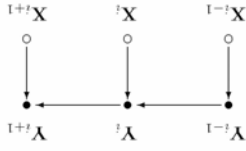


Figure 2. Graphical structures of simple HMMs (left), MEMMs (center), and the chain-structured case of CRFs (right) for sequences. An open circle indicates that the variable is not generated by the model.

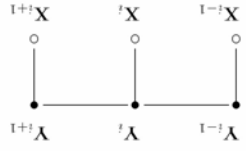
HMM



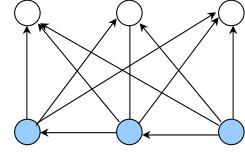
MEMM



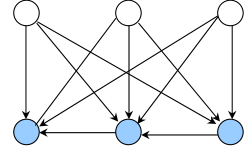
CRF



NO!



YES



YES

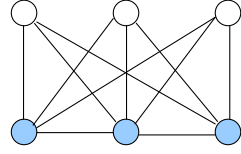


Figure 2. Graphical structures of simple HMMs (left), MEMMs (center), and the chain-structured case of CRFs (right) for sequences. An open circle indicates that the variable is not generated by the model.

Maximum Entropy II

- Also: regularization (smoothing)
- Maximize likelihood = Minimize "log-loss"

$$\min_k \|w\|_2^2 - \sum_i \left(w^\top f_i(y_i) - \log \sum_y \exp(w^\top f_i(y)) \right)$$

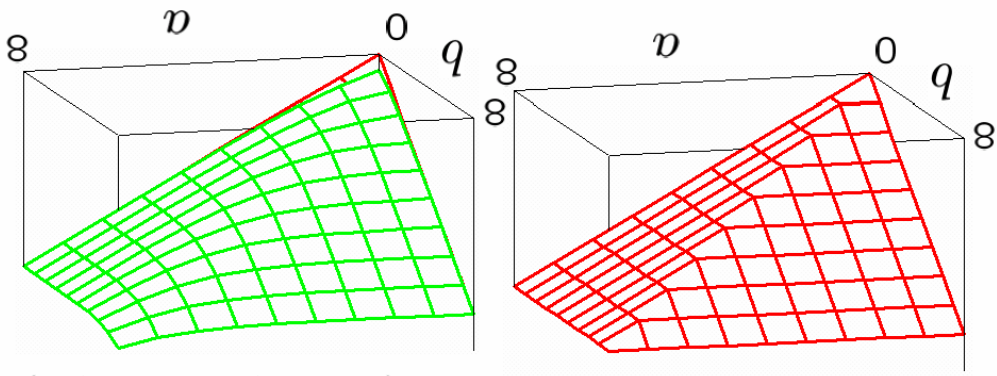
- Motivation
- Connection to maximum entropy principle
- Might want to do a good job of being uncertain on noisy cases...
- ... in practice, though, posteriors are pretty peaked



"Soft-Max"

$$\max(a, b) \approx \log(\exp(a) + \exp(b))$$

$$\max(a, b)$$



Max vs "Soft-Max" Margin

- SVMs:
- Maxent:

$$\min_k \|w\|_2^2 - \sum_i \left(w^\top f_i(y_i) - \log \sum_y \exp(w^\top f_i(y)) \right)$$

Soft Margin

- Very similar! Both try to make the true score better than a function of the other scores.
- The SVM tries to beat the augmented runner-up
- The maxent classifier tries to beat the "soft-max"

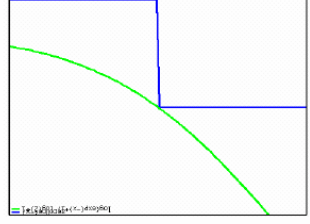


Log-Loss

- If we view maxent as a minimization problem:

$$\min_k \|w\|_2^2 - \sum_i \left(w^\top f_i(y_i) - \log \sum_y \exp(w^\top f_i(y)) \right)$$

- This minimizes the "log-loss" on each example



$$-\log \left(\frac{\sum_y \exp(w^\top f_i(y))}{\exp(w^\top f_i(y_i))} \right) = -\log P(y_i | x_i, w)$$

- Log-loss bounds zero-one loss

$$w^\top f_i(y_i) - \max_{y \neq y_i} w^\top f_i(y)$$

Loss Functions: I



- Zero-One Loss

$$\sum_{step} \left(w^T f_i(y_i) - \max_{y \neq y_i} w^T f_i(y) \right)$$

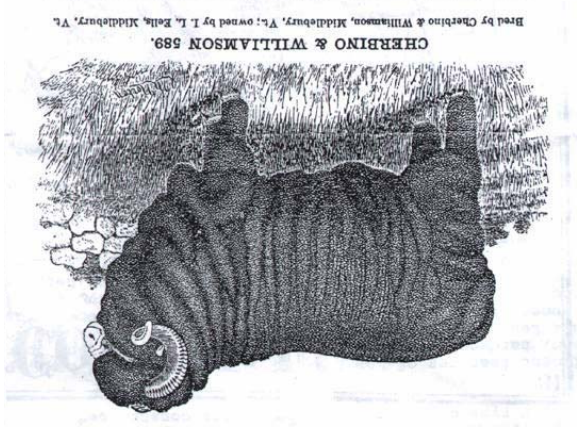
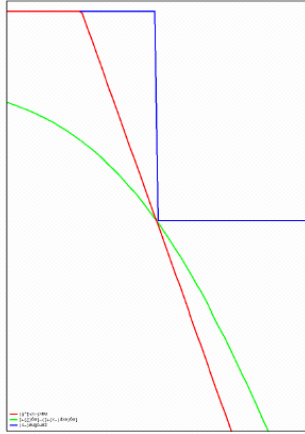
- Hinge

$$\sum_i \left(w^T f_i(y_i) - \max_y [w^T f_i(y) + \ell(y)] \right)$$

- Log

$$\sum_i \left(w^T f_i(y_i) - \log \sum_y \exp(w^T f_i(y)) \right)$$

$$w^T f_i(y_i) - \max_{y \neq y_i} w^T f_i(y)$$



Merino Sheep