# Natural Language Processing
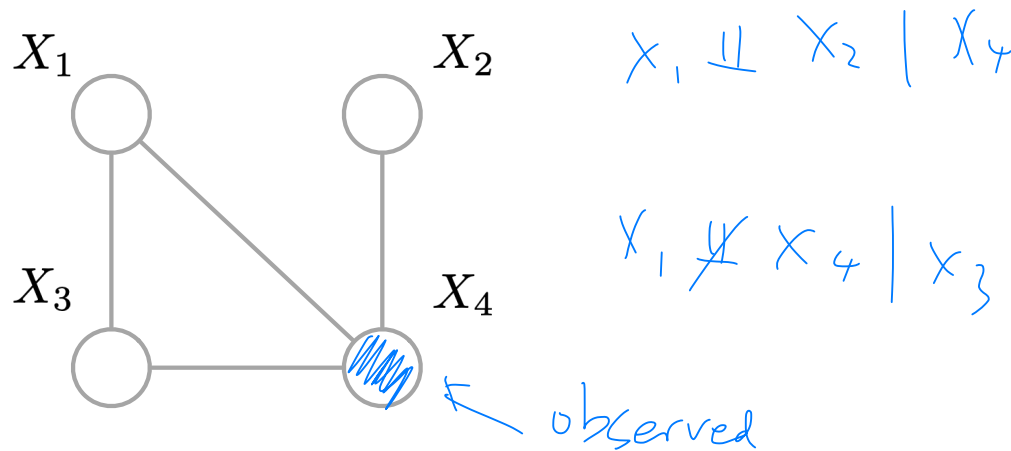## CSE 325/425



Sihong Xie

Lecture 11:
- Conditional random field
- Neural network revisit (forward propagation)
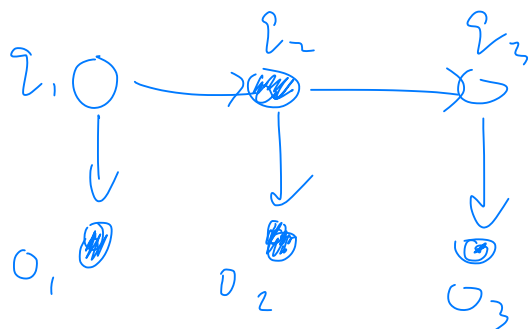
# Graphical models

Conditional independence in a graphical model:

- $A \perp\!\!\!\perp B \mid C$   if any path from A to B have to pass some variables in C.



$X_1$  $X_2$

$X_3$  $X_4$

$X_1 \perp\!\!\!\perp X_2 \mid X_4$

$X_1 \not\perp\!\!\!\perp X_4 \mid X_3$

← observed

"blocks" all paths from $X_1$ to $X_2$

$\sum_{z_1} \sum_{z_3} Pr(z_1, z_3 \mid z_2) = x$      $\left(\sum_{z_1} Pr(z_1 \mid z_2)\right)$

$z_1 \perp\!\!\!\perp z_3 \mid z_2 \implies Pr(z_1, z_3 \mid z_2)$      $\left(\sum_{z_3} Pr(z_3 \mid z_2)\right)$

"forward / backward" $= Pr(z_1 \mid z_2) \times Pr(z_3 \mid z_2)$

Hmm:

$z_1 \bigcirc \longrightarrow \underset{z_2}{\bullet} \longrightarrow \underset{z_3}{\bigcirc}$

$O_1$   $O_2$   $O_3$
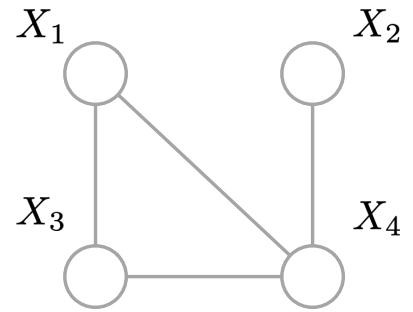
Example of cliques:

① clique = $\{X_1, X_3, X_4\}$

(maximal clique)

② clique = $\{X_2, X_4\}$    not a maximal clique

③ clique = $\{X_1, X_3\}$ ← not a maximal clique

$X_1 \perp\!\!\!\perp X_2 | X_4$

# Cliques

Cliques and maximal cliques

- A clique **c** is a set of nodes that are fully connected
  - ∩ any two nodes in the clique are connected by an edge.
  - '' two nodes not in a clique **can** become conditional independent.

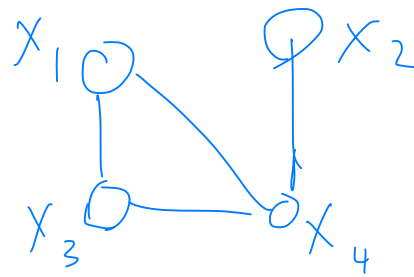$$\Pr(X_i, X_j | X_{\setminus i,j}) = \Pr(X_i | X_{\setminus i,j})\Pr(X_j | X_{\setminus i,j})$$

- A maximal clique **c** is a clique that adding any new node

  will make their union not a clique.
  - ∘ the collection of maximal cliques on a graphical model encode all

    conditional independence properties.

when can do

"polynomial time"

marginalization,

such as,

$Pr(X_1)$

$X_1$        $X_2$

$X_3$        $X_4$

$$\sum_{X_2 X_3 X_4} \Pr(X_1 X_3)\Pr(X_1 X_4) \; P(X_3 X_4)\Pr(X_2 X_4)$$

$$= \sum_{X_2 X_3 X_4} \Pr(X_1 X_2 X_3 X_4)$$

$$= \sum_{X_2 X_3 X_4} \Pr(X_1 X_3 X_4) \times \Pr(X_2 X_4)$$

$X_1$ $X_2$
$X_3$ $X_4$

$c = 1 : \{X_1, X_3, X_4\}$

$c = 2 : \{X_2, X_4\}$

$\mathcal{C} = \{1, 2\}$

# Factorization

Factorization using cliques.

$X_c = \{x_i : X_i \in c\}$

$$\Pr(X_1, \ldots, X_n) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \psi_c(X_c)$$

$\psi_1(X_1, X_3, X_4) : (x_1, x_3, x_4) \rightarrow (0, \infty)$

$\psi_2(X_2, X_4) : (x_2, x_4) \rightarrow (0, \infty)$

where

$\mathcal{C} = $ set of all cliques

- **c** is a clique.

- $\psi_c(X_c) \geq 0$ is a potential function of the variable in the clique **c**.

  - this is not a joint distribution of the variables $X_c$

partition function
- The normalization factor is defined as

$O(2^4)$

$$Z = \sum_{x_1, \ldots, x_n} \prod_{c \in \mathcal{C}} \psi_c(X_c)$$

$X_1, X_2, Y_3, \& X_4 \in \{0, 1\}$

$Z = \psi_1(X_1 = 0, X_3 = 0, X_4 = 0) \times \psi_2(X_2 = 0, X_4 = 0)$

$+ \psi_1(X_1 = 1, X_3 = 0, X_4 = 0) \times \psi_2(X_2 = 0, X_4 = 0)$

$\psi_1(X_1 = X_3 = X_4 = 1)$

$\psi(X_2 = X_4 = 1)$

# Conditional random fields

Conditional random fields:

- A "random field" refer to "a set of dependent random variables".
  - a specific type of "graphical models".

*observed words from a sentence*

- "Conditional" means "conditioning on observed data"
  - making CRF a discriminative model (vs. generative models such as HMM).

- Use the maximum entropy principle – a generalization of logistic regression.
  - each factor is in the form
    - **c** is a maximum clique
$$\psi_c(X_c; O) = \exp\left\{\sum_{i=1}^{d} \theta_i f_i(X_c; O)\right\}$$

*CRF parameter*
$$\vec{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

  - joint conditional distribution

$$\Pr(X_1, \ldots, X_n; O) = \frac{1}{Z(O)} \prod_{c \in \mathcal{C}} \exp\left\{\sum_{i=1}^{d} \theta_i f_i(X_c; O)\right\}$$

$q_1 \cdots q_T$    $Q_c$    $\in \mathbb{R}^d$

$q_t \in \{1, \ldots, N\}$

Summation takes

$Z(O) = \sum_{q_1 \cdots q_T} \prod_{c \in \mathcal{C}} \exp\left\{\sum_{i=1}^{d} \theta_i f_i(Q_c; O)\right\}$
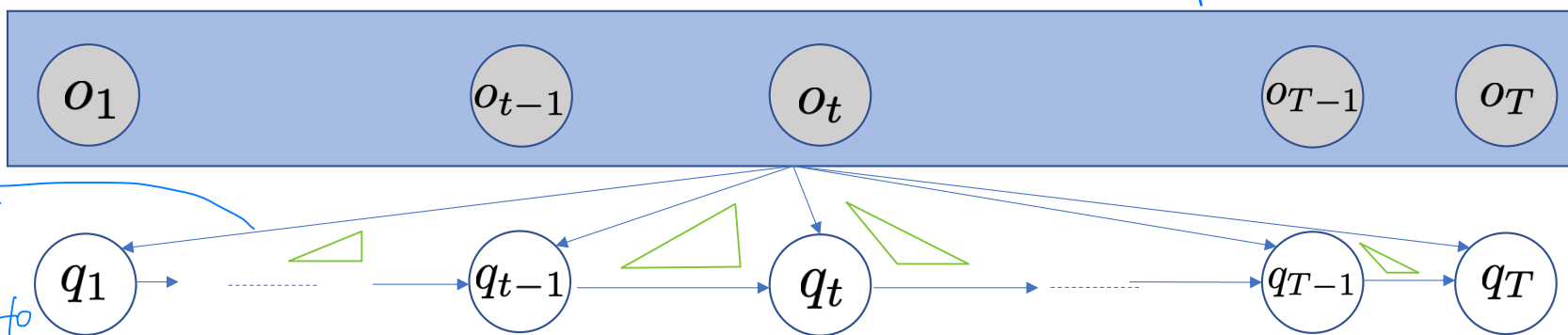
$O(N^T)$

$(q_t, q_{t-1})$

# CRF for POS tagging

The graphical model is a linear chain

want to
predict
each tag
using info
of the
entire
sentence



All factors (or equivalently, maximum cliques) are pairwise

$$\psi_{t-1,t}(q_t, q_{t-1}; O) = \exp\left\{\sum_{i=1}^d \theta_i \boxed{f_i(q_t, q_{t-1}; O)}\right\}$$

$\underbrace{\phantom{\psi_{t-1,t}}}_{Q_+}$

$: (q_t, q_{t-1})$
$\rightarrow (0, \infty)$

$\mathcal{C} = \left\{ \{\xi_1, \xi_2, O\} \right.$

$\{\xi_2, q_3, O\},$

The inner product = the compatibility score of the two tags:

$$s_t(q_{t-1}, q_t; \boldsymbol{\theta}, O) = \langle \boldsymbol{\theta}, \mathbf{f}(q_{t-1}, q_t) \rangle \in \mathbb{R}$$
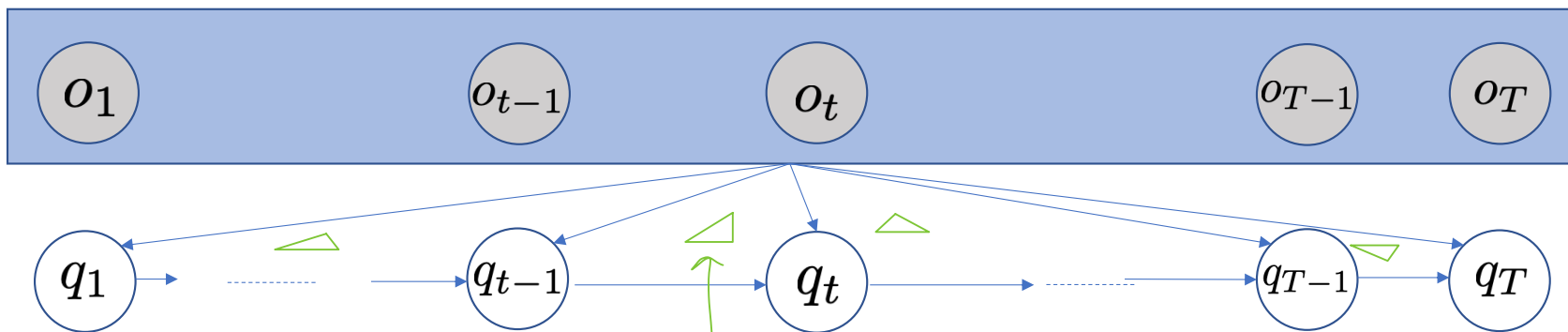
(replacing

$a_{\xi_{t-1}\xi_t} b_{\xi_t}(o_t)$

in HMM, or

$Pr(\xi_t | q_{t-1}, O)$ in MEMM)

$\vec{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \in \mathbb{R}^d,$

$\vec{\hat{f}}$

$\mathbf{f}(\xi_t, \xi_{t-1}, O) = \begin{bmatrix} f_1(q_t, q_{t-1}, O) \\ \vdots \\ f_d(q_t, q_{t-1}, O) \end{bmatrix}$

$\left. \{\xi_{T-1}, \xi_T, O\} \right\}$

# CRF for POS tagging

The graphical model is a linear chain



- The joint distribution of a tag sequence, conditioned on a word sequence is

$$\text{Pr}(q_1, \ldots, q_T; O) = \frac{1}{Z(O)} \prod_{t=1}^{T} \exp\left\{ \sum_{i=1}^{d} \theta_i f_i(q_t, q_{t-1}; O) \right\}$$

decompose over steps of the sequence:
=> conditional independence
=> polynomial inference alg.

$$= \frac{1}{Z(O)} \exp\left\{ \sum_{t=1}^{T} \sum_{i=1}^{d} \theta_i f_i(q_t, q_{t-1}; O) \right\} = \frac{1}{Z(O)} \exp\left\{ \sum_{t=1}^{T} s_t(q_{t-1}, q_t; \boldsymbol{\theta}, O) \right\}$$

$$\langle \vec{\theta}, \vec{f}(q_t, q_{t-1}; O) \rangle$$

# Predicting sequence using CRF

Input:

- an input sentence $O = [o_1, \ldots, o_T], o_t \in V$

- and a trained CRF model $\boldsymbol{\theta}$

Output:

- optimal POS tag sequence $Q^* = \underset{Q}{\arg\max} \Pr(Q|O; \boldsymbol{\theta})$ $= \underset{Q}{\arg\max} \frac{1}{Z(o)} \exp\left\{ \sum_{t=1}^{T} s_t(\Sigma_{t-1}, \Sigma_{t}; \theta, o) \right\}$

$$= \underset{Q}{\arg\max} \sum_{t=1}^{T} s_t(q_{t-1}, q_t; \boldsymbol{\theta}, O)$$

Adapt the Viterbi algorithm for HMM to CRF prediction:

- change the scores in HMM $\quad s_t(q_{t-1} = i, q_t = j; \lambda, O) = a_{i,j} b_j(o_t)$

to the scores defined for CRF.

? MEMM    what you're replacing

$$\log Z(o) = \log \sum_Q \exp\left\{ \sum_{t=1}^{T} \langle \vec{\theta}, \vec{f}(q_t, q_{t-1}, \theta, o) \rangle \right\}$$

$$\frac{\partial}{\partial \theta} \log Z(o) = \frac{1}{Z(o)} \frac{\partial}{\partial \theta} Z(o) = \frac{1}{Z(o)} \sum_Q \exp\left\{ \sum_{t=1}^{T} \langle \vec{\theta}, \vec{f} \rangle \right\} \sum_{t=1}^{T} \frac{\partial}{\partial \theta} \langle \vec{\theta}, \vec{f}(q_t, q_{t-1}, \theta, o) \rangle$$

# Learning a CRF model

$$= \sum_Q Pr(Q|\theta, o) \sum_{t=1}^{T} \vec{f}(q_t, q_{t-1}, \theta, o)$$

$$= \sum_{t=1}^{T} \mathbb{E}_{Q \sim Pr(Q|\theta, o)} \vec{f}(q_t, q_{t-1}, \theta, o)$$

Not much more difficult than training a logistic regression model!

- Input: $m$ POS-tagged sentences.

*Expectation of feature vector* $\mathbb{E}_{\vec{x} \sim Pr(\vec{x})} \vec{f}(\vec{x})$
$= \sum_{\vec{x}} Pr(\vec{x}) \vec{f}(\vec{x})$

- MLE:
$$\boldsymbol{\theta}^* = \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{m} \log \Pr(Q^{(i)}|O^{(i)}, \boldsymbol{\theta}) = \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{m} \ell(Q^{(i)}|O^{(i)}, \boldsymbol{\theta})$$

- There is no closed form solution for the parameter, and gradient descent is needed.

*Logistic regression*

$$\frac{\partial}{\partial \boldsymbol{\theta}} \ell(Q|O, \boldsymbol{\theta}) = \frac{\partial}{\partial \theta} \log Pr(Q|O; \theta) = \frac{\partial}{\partial \theta} \log \frac{1}{Z(o)} \exp\left\{ \sum_{t=1}^{T} \langle \vec{\theta}, \vec{f}_t \rangle \right\}$$

*"Glove"*

$$= \frac{\partial}{\partial \theta} \left[ \sum_{t=1}^{T} \langle \vec{\theta}, \vec{f}_t \rangle - \log Z(o) \right]$$

$$= \sum_{t=1}^{T} \vec{f}_t - \sum_{t=1}^{T} \mathbb{E}_{Q \sim Pr(Q|O,\theta)} \vec{f}(q_t, q_{t-1}, \theta, o)$$

decompose over steps of the sequence:
=> conditional independence
=> polynomial inference alg.

*matrix calculus:*

- Recall the gradient for multi-class logistic regression …

*learning rate 0.1*

$$\frac{\partial}{\partial \theta} \langle \vec{\theta}, \vec{f} \rangle = \vec{f}$$

$$\vec{\theta} \leftarrow \vec{\theta} + \eta \frac{\partial}{\partial \theta} \ell(Q|O, \theta) = \vec{\theta} + \eta \left[ \sum_{t=1}^{T} \vec{f}_t - \sum_{t=1}^{T} \mathbb{E}_{Q \sim Pr} \vec{f} \right]$$

# Learning a CRF model

A running example (Cheating Casino)

$$\vec{f} = \begin{bmatrix} f_1 \\ \vdots \\ f_d \end{bmatrix} = \begin{bmatrix} f_1 \begin{pmatrix} \text{no cheating} \\ \text{cheat} \end{pmatrix}; H,T,H \\ = 1 \\ f_2 \begin{pmatrix} \text{cheating} ; H,T,H \\ \text{cheating} \end{pmatrix} \\ \vdots \end{bmatrix}$$
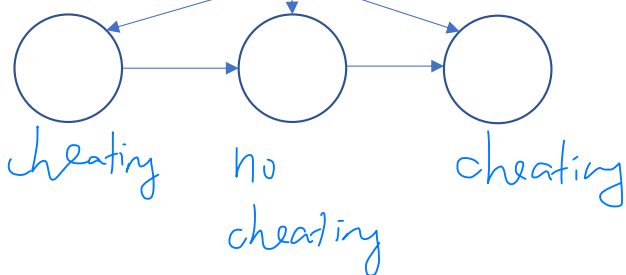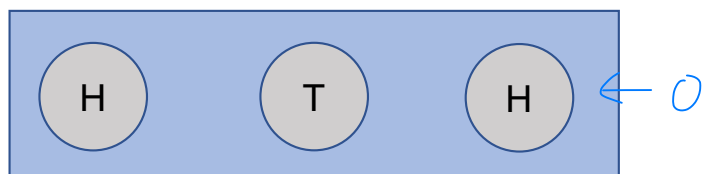
$q_{t-1}$

$q$

$q_t$

Alg: Random $\vec{\theta} = [0.1, 0.1, 0.1 \cdots 0.1] \in \mathbb{R}^d$
Init

# gradient ascent:

for $i = 1 \cdots n$

$$\vec{\theta}^{(i+1)} \leftarrow \vec{\theta}^{(i)} + g\left( \sum_{t=1}^{3} \vec{f}(q_{t-1}, q_t, 0) - \sum_{t=1}^{3} \mathbb{E}_{\substack{Q \\ \sim Pr(Q|0,\theta)}} \vec{f}(q_{t-1}, q_t, 0) \right)$$

flips
of
coins



cheating    no          cheating
            cheating

$q_1$        $q_2$        $q_3$
$t=1$       $t=2$       $t=3$

$$\mathbb{E}_{Q \sim Pr(Q|0,\theta)} \vec{f}(q_1, q_2, 0)$$

$$= \sum_{Q = [q_1, q_2, q_3]} Pr(q_1, q_2, q_3 | 0, \theta) \vec{f}(q_1, q_2, 0)$$

$$= \sum_{q_3} Pr(q_3 | q_1, q_2, 0, \theta) \underbrace{\sum_{q_1, q_2} Pr(q_1, q_2 | 0, \theta)}_{=1} \vec{f}(q_1, q_2, 0)$$

what is $Pr(q_1, q_2 | 0, \theta)$ ? see the next page.

$$\Pr(q_1, q_2; \theta, 0)$$

$$= \sum_{q_3} \frac{1}{Z(0)} \exp\left\{ \sum_{t=1}^{3} S_t(q_{t-1}, q_t; \theta, 0) \right\}$$

↙ dummy fixed POS tag.

$$= \sum_{q_3} \frac{1}{Z(0)} \left[ \exp\{ S_1(q_0, q_1; \theta, 0) \} \right.$$

$$\times \exp\{ S_2(q_1, q_2; \theta, 0) \}$$

$$\left. \times \exp\{ S_3(q_2, q_3; \theta, 0) \} \right.$$

$$= \frac{1}{Z(0)} \underline{\frac{\exp\{ S_1(q_0, q_1; \theta, 0) \}}{}} \quad \rightarrow \text{forward Prob.}$$
$$= \alpha_1(q_1)$$

$$\times \exp\{ S_2(q_1, q_2; \theta, 0) \}$$

$$\times \underbrace{\sum_{q_3} \exp\{ S_3(q_2, q_3; \theta, 0) \}}_{}$$

backward prob. $= \beta_2(q_2)$

$$= \frac{1}{Z(0)} \alpha_1(q_1) \exp\{ S_2(q_1, q_2; \theta, 0) \} \beta_2(q_2)$$

$$Z(0) = \sum_{q_1 q_2 q_3} \exp\{ S_1 + S_2 + S_3 \}$$

$$= \sum_{q_1 q_2 q_3} \exp\{ S_1(q_0, q_1) + S_2(q_1, q_2) + S_3(q_2, q_3) \}$$

$$= \alpha_2(q_2)$$

$$= \sum_{q_3} \left[ \left[ \sum_{q_2} \underbrace{\sum_{q_1} \exp\{ S_1 + S_2 \}}_{} \right] \exp\{ S_3(q_2, q_3) \} \right]$$

$\alpha_3(q_3)$ ←