

# Conditional random field (条件随机场)

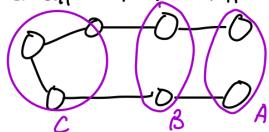
## 定向图 Review

### 三. 无向图 - 马尔科夫网络 Markov Network (Markov Random Field)

一. 条件独立性  $MRF \subseteq \text{Gibbs Distribution}$

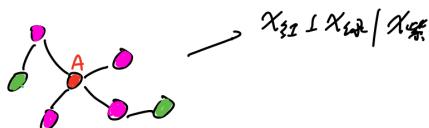
条件独立性可以体现在以下三个方面:

① Global Markov  $x_A \perp x_C \mid x_B$ .  $P(x_A, x_C \mid x_B) = P(x_A \mid x_B) \cdot P(x_C \mid x_B)$



② Local Markov

$x_A \perp x_D \mid x_B$ , 其中 A 是变量, B 是其邻居, D 为其余变量.



③ Pairwise Markov

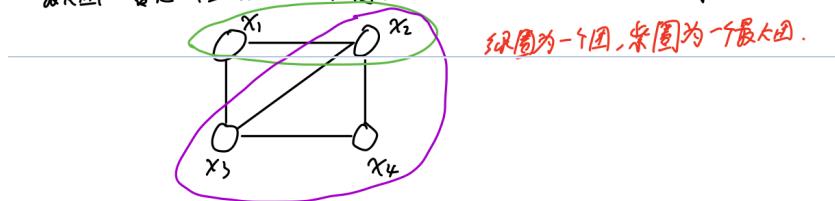
$x_A \perp x_B \mid x_{-\{A,B\}}$ , 即给定所有其他变量两个邻接变量条件独立.

### 二. 团与分解

① 基本概念. 团

团: 一个关于结点的集合, 该集合的结点之间相互连通

最大团: 若在一个团中加入其他结点则不再连通, 则称该团为最大团.



每对相关结点之间定义势函数, 在 Markov Random Field 中, 多个变量的联合分布能表示为团分解为多个势函数的乘积, 每个团对应一个势函数.

$$P(x) = \frac{1}{Z} \prod_{i=1}^k \psi_i(x_{c_i})$$

$C_i$ : 最大团

Explanation:

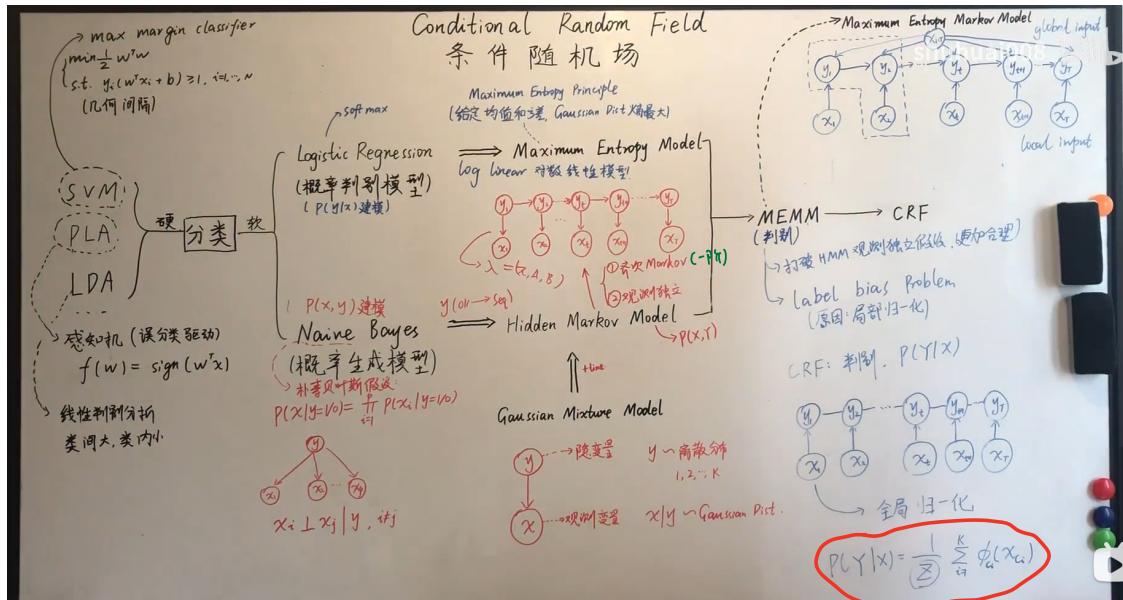
$x_{c_i}$ : 最大团随机变量集合

$$\sum_i P(x) = 1$$

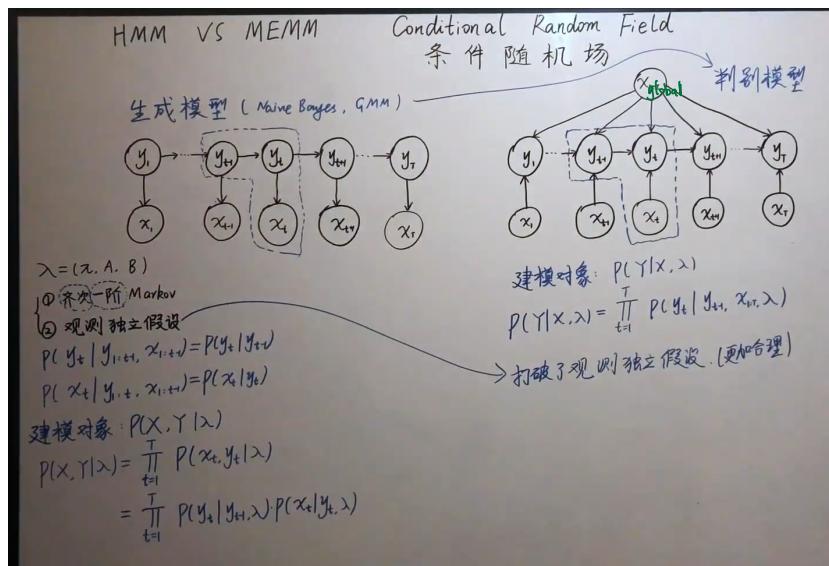
$\psi_i(x_{c_i})$ : 势函数, 必须为正

Z: 归一化因子.  $Z = \sum_i \prod_{j=1}^k \psi_j(x_{c_j})$

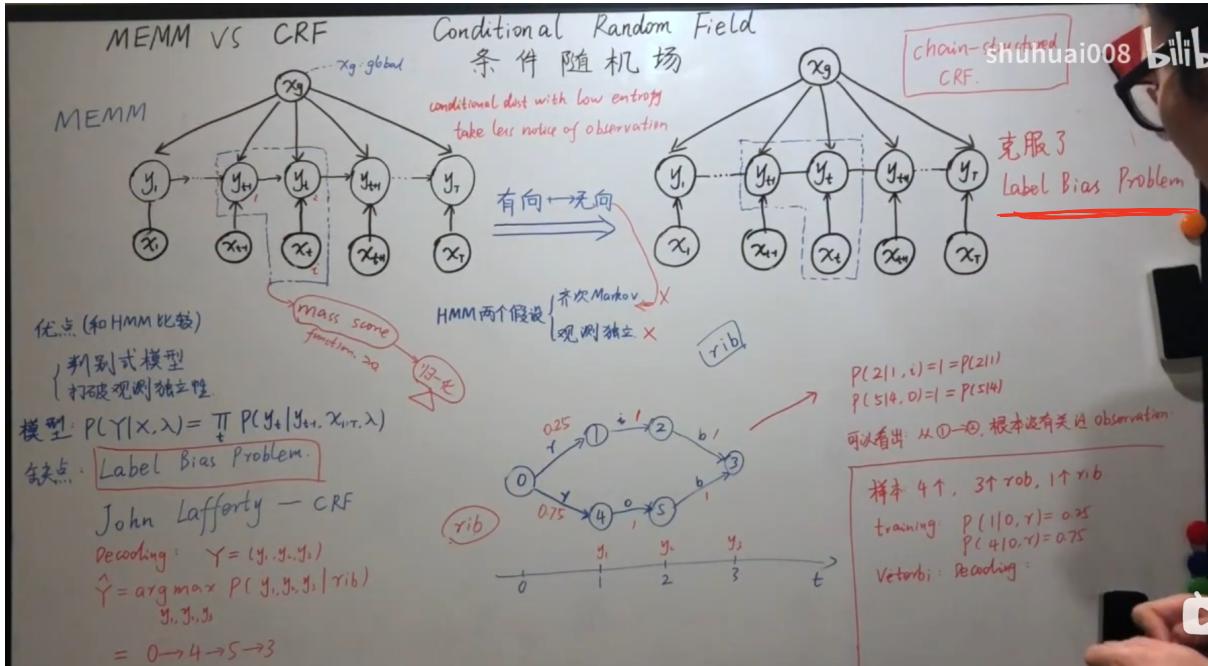
## Background



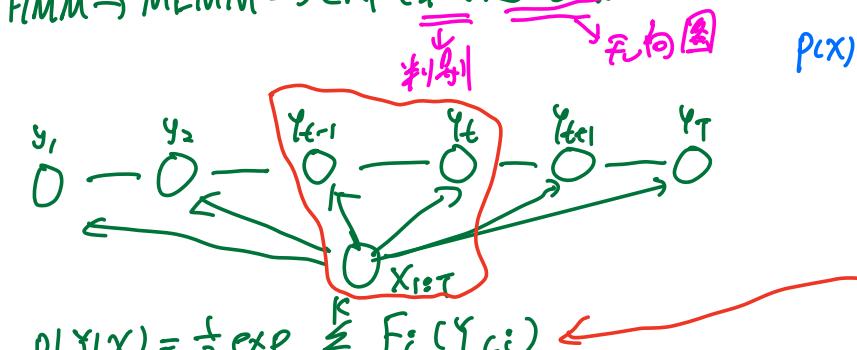
## HMM VS MEMM (maximum-entropy Markov model)



## MEMM VS CRF



HMM → MEMM → CRF (条件随机场)



$$\begin{aligned}
 F(y_{t-1}, y_t, x_{1:T}) &= \Delta y_{t-1}, x_{1:T} + \Delta y_t, x_{1:T} + \Delta y_{t+1}, y_t, x_{1:T} \\
 &= \Delta y_t, x_{1:T} + \Delta y_{t+1}, y_t, x_{1:T}
 \end{aligned}$$

$$\Delta y_{t-1}, y_t, x_{1:T} = \sum_{k=1}^K \lambda_k f_k(y_{t-1}, y_t, x_{1:T})$$

$$\text{同理 } \Delta y_t, x_{1:T} = \sum_{l=1}^L \eta_l g_l(y_t, x_{1:T})$$

Potential functions:

$$p(x) = \frac{1}{Z} \sum_{i=1}^K \psi_i(x_{C_i})$$

$$= \frac{1}{Z} \sum_{i=1}^K \exp[-E(x_{C_i})]$$

$$= \frac{1}{Z} \exp \sum_{i=1}^K F_i(x_{C_i})$$

$C_i$ : 最大团

$K$ :  $K$  个最大团

$f_k, g_l$ : 相应的特征函数

$\lambda_k, \eta_l$ : 系数

$$\text{P}(y|x) = \frac{1}{Z} \exp \left[ \sum_{t=1}^T \lambda_k f_k(y_{t-1}, y_t, x_{1:T}) + \sum_{t=1}^L \eta_t g_t(y_t, x_{1:T}) \right]$$

CRF 的 PDF \*

试图化简上式 PDF，用向量化表示：

首先想去掉内部的  $\sum_{k=1}^K$  和  $\sum_{t=1}^L$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_T \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_T \end{bmatrix} \quad \lambda = \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_K \end{bmatrix} \quad \eta = \begin{bmatrix} \eta_1 \\ \eta_2 \\ \vdots \\ \eta_L \end{bmatrix}$$

$$f = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_K \end{bmatrix} = f(y_{t-1}, y_t, x_{1:T}) \quad g = \begin{bmatrix} g_1 \\ g_2 \\ \vdots \\ g_K \end{bmatrix} = g(y_t, x_{1:T})$$

$$\Rightarrow P(y|x) = \frac{1}{Z} \exp \left[ \sum_{t=1}^T \left[ \lambda^T f(y_{t-1}, y_t, x_{1:T}) + \eta^T g(y_t, x_{1:T}) \right] \right]$$

然后进一步想去掉  $\sum_{t=1}^T$ , 因为  $f, g$  项有关：

$$P(y|x) \approx \frac{1}{Z} \exp \left[ \lambda^T \sum_{t=1}^T f(y_{t-1}, y_t, x_{1:T}) + \eta^T \sum_{t=1}^T g(y_t, x_{1:T}) \right]$$

$$\Theta = \begin{bmatrix} \lambda \\ \eta \end{bmatrix}_{K \times L} \quad H = \begin{bmatrix} \sum_{t=1}^T f(y_{t-1}, y_t, x_{1:T}) \\ \sum_{t=1}^T g(y_t, x_{1:T}) \end{bmatrix}_{K \times L}$$

$$\therefore P(y|x) = \frac{1}{Z} \exp \left( \Theta^T H(y_{t-1}, y_t, x_{1:T}) \right)$$

## 模型要解决的问题

Learning: 参数学习出来 (parameter estimation)  $\Rightarrow \hat{\theta} = \arg \max_{\theta} \sum_{i=1}^N p(y_i | x_i)$   
 Inference:  $\left\{ \begin{array}{l} \text{marginal prob: } p(y_t | X) \rightarrow p(y_t = i | X, \theta) \\ \text{conditional prob: 生成模型的后验. 这里没有} \\ \text{MAP inference: decoding: } \hat{y} = \arg \max_{y=y_1, y_2, \dots, y_T} p(y | x) \end{array} \right.$

## Review:

### 1. Variable inference

有向图:



$$P(E) = \sum_{A,B,C,D} P(A, B, C, D, E)$$

逆推:

$$P(E) = \sum_{A,B,C,D} P(A) \cdot P(B|A) \cdot P(C|B) \cdot P(D|C) \cdot P(E|D)$$

求 A:

$$= \sum_{B,C,D} P(C|B) P(D|C) P(E|D) \frac{\sum_A P(A) P(B|A)}{\sum_A P(A)}$$

$$= \sum_{B,C,D} P(C|B) P(D|C) P(E|D) P(B)$$

求 B:

$$= \sum_{C,D} P(D|C) P(E|D) \frac{\sum_B P(B) P(C|B)}{\sum_B P(B)}$$

$$= \sum_{C,D} P(D|C) P(E|D) P(C)$$

求 C:

$$= \sum_D P(E|D) \sum_C P(C) P(D|C)$$

$$= \sum_D P(E|D) P(D)$$

$O(kn^2)$   $k$  is the number of variables

Forward:



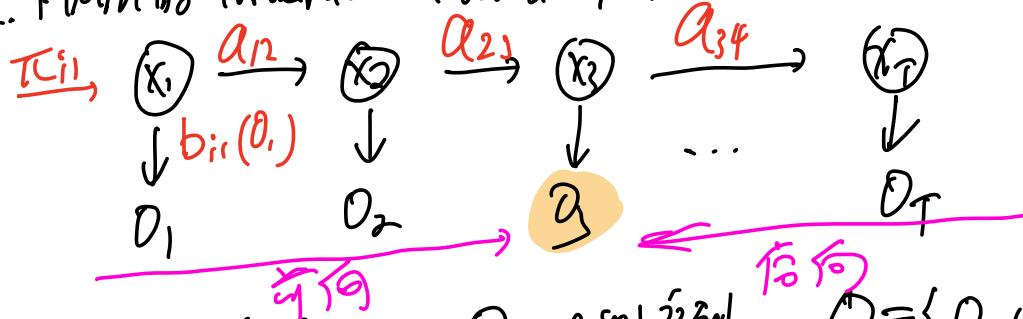
$$p(x) = \frac{1}{Z} \prod_{i=1}^k \psi_i(x_{c_i})$$

$$\begin{aligned} P(e) &\propto \sum_{a,b,c,d} \phi(a,b)\phi(b,c)\phi(c,d)\phi(d,e) \\ &= \sum_{b,c,d} \phi(b,c)\phi(c,d)\phi(d,e) \sum_a \phi(a,b) \\ &= \sum_{b,c,d} \phi(b,c)\phi(c,d)\phi(d,e)m_a(b) \\ &= \sum_{c,d} \phi(c,d)\phi(d,e) \sum_b \phi(b,c)m_a(b) \\ &= \sum_{c,d} \phi(c,d)\phi(d,e)m_b(c) \\ &= \sum_d \phi(d,e) \sum_c \phi(c,d)m_b(c) \\ &= \sum_d \phi(d,e)m_c(d) \\ &= \sum_d \phi(d,e)m_c(d) \\ &= m_d(e) \end{aligned}$$

Finally we normalize to obtain a proper probability:

$$P(e) = \frac{m_d(e)}{\sum_e m_d(e)}$$

2. HMM has forward-backward algorithm



Target:  $P(O|\lambda)$

$$O: \underline{o_1, o_2, o_3, \dots, o_T}$$

$$I: \underline{i_1, i_2, i_3, \dots, i_T}$$

$$O = \{o_1, o_2, o_3, \dots, o_T\}$$

$$I = \{i_1, i_2, i_3, \dots, i_T\}$$

$$P(I|\lambda) = \pi_{i_1} \alpha_{i_1 i_2} \alpha_{i_2 i_3} \cdots \alpha_{i_{T-1} i_T}$$

$$p(O|I, \lambda) = b_{i1}(O_1) b_{i2}(O_2) b_{i3}(O_3) \cdots b_{iT}(O_T)$$

$$p(O, I|\lambda) = p(I|\lambda) \cdot p(O|I, \lambda)$$

$$= \prod_{i1} a_{i1, i2} b_{i1}(O_1) \cdots \prod_{iT} a_{iq, iT} b_{iT}(O_T)$$

$$p(O|\lambda) = \sum_I \prod_{i1} a_{i1, i2} b_{i1}(O_1) \cdots \prod_{iT} a_{iq, iT} b_{iT}(O_T)$$

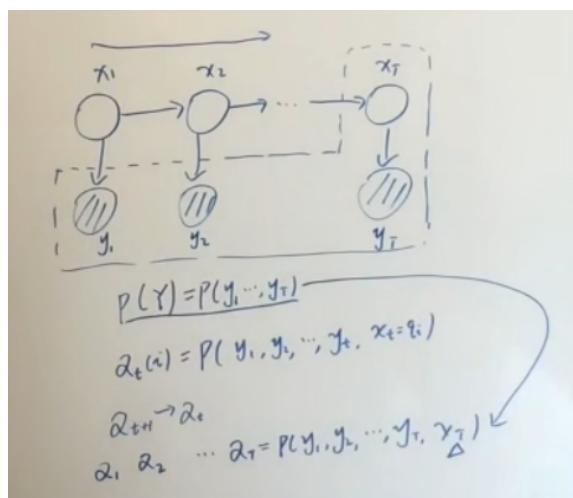
$$\boxed{p(O|\lambda) = \sum_{i1, i2, i3, \dots, iT} \prod_{i1} a_{i1, i2} b_{i1}(O_1) \cdots \prod_{iT} a_{iq, iT} b_{iT}(O_T)}$$

$\Rightarrow x: a_t(i) = p(O_1, O_2, \dots, O_T, i_t = q_i|\lambda)$

VE~~法~~:

$$p(O|\lambda) = \sum_{i2 \dots iT} \cdots \sum_{i1} \prod_{i1} a_{i1} \underline{a_{i2}} b_{i1}(O_1)$$

$$a_{t+1}(i) = \sum_{j=1}^N a_t(j) a_{ji} b_i(O_{t+1}) \quad i=1, 2, \dots, N$$



Learning: Given  $P(Y=y|X=x)$ , find  $P(y_t=i|X)$

$$\begin{aligned} P(y_t=i|X) &= \sum_{y_1, y_2, \dots, y_{t-1}, y_t, y_T} P(Y(X)) \\ &= \sum_{y_1, \dots, y_{t-1}} \sum_{y_t \in S} \frac{1}{2} \prod_{t'=1}^T \psi_t(y_{t-1}, y_t, x) \\ &= \Delta_L \Delta_R \end{aligned}$$

其中  $\Delta_L$ :  $\sum_{t=1}^T \psi_t(y_t, y_{t-1}, x) \dots \dots \sum_{y_1} \psi_1(y_1, y_0, x) \leq \underbrace{\psi_1(y_1, y_0, x)}_{y_0} \leq \underbrace{\psi_r(y_r, y_0, x)}_{y_0}$

$$\text{记 } \Delta_L = \alpha_t(i) = \sum_{j \in S} \psi_t(y_{t-1}=j, y_t=i, x) \alpha_{t-1}(j)$$

Similarly, 得到

$$\Delta_R = \beta_t(i) = \sum_{j \in S} \psi_{t+1}(y_t=i, y_{t+1}=j, x) \beta_{t+1}(j)$$

本质: VE 法

从  $\alpha$  到 forward-backward 法 !

## Learning

target:  $\hat{\theta} = \arg \max_{\theta} \prod_{i=1}^N p(y^{(i)} | x^{(i)}) \quad \theta = (\lambda, \eta)$

在进行各种类型的推断之前，还需要对参数进行学习：

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \prod_{i=1}^N p(y^i | x^i) \\ &= \arg \max_{\theta} \sum_{i=1}^N \log p(y^i | x^i) \\ &= \arg \max_{\theta} \sum_{i=1}^N [-\log Z(x^i, \lambda, \eta) + \sum_{t=1}^T [\lambda^T f(y_{t-1}, y_t, x) + \eta^T g(y_t, x)]]\end{aligned}$$

上面的式子中，第一项是对数配分函数，根据指数族分布的结论：

$$\nabla_{\lambda} (\log Z(x^i, \lambda, \eta)) = \mathbb{E}_{p(y^i | x^i)} [\sum_{t=1}^T f(y_{t-1}, y_t, x^i)]$$

其中，和  $\eta$  相关的项相当于一个常数。求解这个期望值：

$$\mathbb{E}_{p(y^i | x^i)} [\sum_{t=1}^T f(y_{t-1}, y_t, x^i)] = \sum_y p(y | x^i) \sum_{t=1}^T f(y_{t-1}, y_t, x^i)$$

第一个求和号的复杂度为  $O(S^T)$ ，重新排列求和符号：

$$\begin{aligned}\mathbb{E}_{p(y^i | x^i)} [\sum_{t=1}^T f(y_{t-1}, y_t, x^i)] &= \sum_{t=1}^T \sum_{y_{1:t-2}} \sum_{y_{t-1}} \sum_{y_t} \sum_{y_{t+1:T}} p(y | x^i) f(y_{t-1}, y_t, x^i) \\ &= \sum_{t=1}^T \sum_{y_{t-1}} \sum_{y_t} p(y_{t-1}, y_t | x^i) f(y_{t-1}, y_t, x^i)\end{aligned}$$

和上面的边缘概率类似，也可以通过前向后向算法得到上面式子中的边缘概率。

于是：

$$\underline{\nabla_{\lambda} L} = \sum_{i=1}^N \sum_{t=1}^T [f(y_{t-1}, y_t, x^i) - \sum_{y_{t-1}} \sum_{y_t} p(y_{t-1}, y_t | x^i) f(y_{t-1}, y_t, x^i)]$$

利用梯度上升算法可以求解。对于  $\eta$  也是类似的过程。

$$\begin{aligned}\lambda^{(t+1)} &= \lambda^{(t)} + \text{step} \cdot \nabla_{\lambda} L(\lambda^{(t)}, \eta^{(t)}) \\ \eta^{(t+1)} &= \eta^{(t)} + \text{step} \cdot \nabla_{\eta} L(\lambda^{(t)}, \eta^{(t)})\end{aligned}$$

