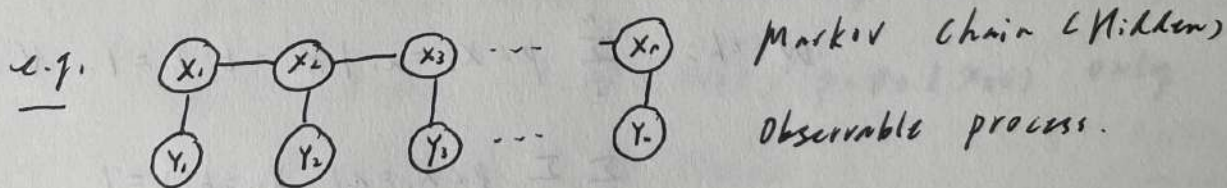# MRFs and HMMs

MRF refers to "Markov Random Field".

HMM refers to "Hidden Markov Model".

Def: i) $(X_s)_{s \in G}$ collection of r.v's is a random field indexed by $G$. Set of nodes of graph.

ii) $s \sim t$ for $s, t \in G$ means they're neighbour.

$$N(t) = \{ s \in G \mid s \sim t \}.$$

iii) $(X_s)_{s \in G}$ is MRF if $p( X_t = x_t \mid X_s = x_s, s \neq t )$

$$= p( X_t = x_t \mid X_s = x_s, s \in N(t) ) \stackrel{\Delta}{=:} p( x_t \mid x_{N(t)} ). \forall t.$$

iv) HMM is a MRF. st. some r.v.'s are observed but others are hidden.

e.g.



Markov Chain (Hidden)

Observable process.

$$p( y_t \mid x, y_{\neq t} ) = p( y_t \mid x_t ).$$

HMMs fit in Bayesian Framework nicely:

For $X$ unknown: $\quad$ prior $= Y \mid X \xrightarrow[\text{Formula}]{\text{Bayesian}}$ Posterior $: X \mid Y$.

We will make a reseasonable estimate for $p(X \mid Y)$

from $Y \mid X$.

(1) <u>Gibbs Dist.</u>:

Next. we will specify a MRF. by 2 methods:

i) <u>Set of Conditional Dist.</u>

$( p(X_t | X_{N(t)}) )_{t \in G}$ may not be a consistent
prob. dist. when we give MRF one dist.

Thm. If we have specified condition dist.
$\mathcal{L}(X_1 | X_2)$ for r.v. $X_1$, $X_2$. $X_1 \in S_1 = (a_i)_i^n$.
$X_2 \in S_2 = (b_i)_i^m$. Then. we're free to
decide one more dist. $\mathcal{L}(X_2 | X_1 = a)$. $a \in S$.

$\underline{pf}$: By Bayesian Formula:

$$p(X_1 = a_i | X_2 = b_j) = p(X_1 = a_i, X_2 = b_j) \Big/ \sum_{k}^{n} p(X_1 = a_k, X_2 = b_j)$$

for $1 \leq i \leq n$, $1 \leq j \leq m$

With: $\sum_{i}^{n} p(X_i = a_i | X_2 = b_j) = 1$. $\forall 1 \leq j \leq m$.

$$\sum_{i} \sum_{j} p(X_1 = a_i, X_2 = b_j) = 1.$$

Consider $( p(X_1 = a_i, X_2 = b_j) )_{i,j}$ as set of
unknown variable. $p(a_i | b_j)$ is known.
The order of linear equation above is
$mn - m + 1$. $\Rightarrow$ At most choose $\mathcal{L}(X_2 | X_1 = a)$

ii) <u>Hammersley - Clifford Thm</u>:

Def: i) Set of nodes $G$ is complete if every distinct nodes are neighbour of each other.

ii) A clique is max set of nodes. st. complete.

iii) $G$ is finite graph. Gibbs dist. w.r.t $G$ is pmf $p(x) = \prod\limits_{\substack{c \text{ is} \\ \text{complete}}} V_c(x)$ & $V_c$ only depends on $X_c = (X_s)_{s \in c}$. for $c$ is clique. $x \in S_g$. (config.)

Rmk: i) $X_c = \eta_c \Rightarrow V_c(x) = V_c(\eta)$.

ii) $p(x)$ can be reduced $= \prod\limits_{c \text{ clique}} V_c(x)$.

Thm. (H-C Thm)

$X = (X_1, X_2 \cdots X_n)$ has positive joint pmf. Then:

$X$ is MRF on $G$ $\iff$ $X$ has a Gibbs dist. on $G$.

Pf: $S_g = S_1 \times S_2 \cdots \times S_n$. $S_k$ is state space of $X_k$.

Denote: $O$ means arbitrary element. (fix)

($\Leftarrow$). Show: $p(X_t \mid X_{\neq t}) / p(O_t \mid X_{\neq t})$ only depends on $X_{N(t)}$.

Note: $p(X_t, X_{\neq t}) / p(O_t, X_{\neq t}) = p(X_t \mid O) / p(O_t \mid O)$

$= \dfrac{\prod_{t \in c} V_c(X_t, X_{\neq t})}{\prod_{t \in c} V_c(O_t, X_{\neq t})} \cdot \dfrac{\prod_{t \notin c} V_c(X_t, X_{\neq t})}{\prod_{t \notin c} V_c(O_t, X_{\neq t})}$

$= \prod_{t \in c} V_c(X_t, X_{\neq t}) / \prod_{t \in c} V_c(O_t, X_{\neq t})$

($\Rightarrow$) We want to write $p(x)$ in form:

$\prod_A V_A(x)$. $V_A \equiv 1$. if $A$ isn't complete.

Set: $p(X_0, O_{p(c)}) = \prod\limits_{A \subseteq D} V_A(x)$. $D \subseteq \{1, 2, \cdots n\}$.

Then we can find $V_a$ reccurrsively.

1) $D = \emptyset$. $p(0) = V_x(x)$

2) $D = \{t\}$. $V_{\{t\}}(x) = p(x_t, 0_{\neq t})/p(0)$

3) $D \subseteq \{1,2,\cdots n\}$. $V_D(x) = p(x_D, 0_{\neq D})/\prod_{\substack{A \subseteq D \\ \neq}} V_A(x)$.

Next. prove: $V_A \equiv 1$. if $A$ not complete.

By induction on $|A|$. $|A| \leq 1$ $\checkmark$.

For $n = k+1$. (Suppose $n \leq k$ holds)

if $t, u \in A$. not neighbour. $A = \{t, u\} \cup B$

Note: $p(x_A, 0_{\neq A}) = p(x_t, x_u, x_B, 0_{\neq A})$

$$= \frac{p(x_t \mid x_u, x_B, 0_{\neq A})}{p(0_t \mid x_u, x_B, 0_{\neq A})} p(0_t, x_u, x_B, 0_{\neq A})$$

$$= \frac{p(x_t \mid x_B, 0_{A^c \cup \{u\}})}{p(0_t \mid x_B, 0_{A^c \cup \{u\}})} p(\square)$$

$$= \frac{\prod_{0 \subseteq B \cup \{t\}} V_D \; \prod_{0 \subseteq B \cup \{t\}} V_D}{\prod_{0 \subseteq B} V_D} = \prod_{\substack{D \subseteq A \\ \neq}} V_D \quad (\text{by induct})$$

prop. $(X, Y)$ is MRF on $G = G_X \cup G_Y$. with neighbour

structure $N_{X \cup Y}$. Then:

i) Marginal dist. of $Y$ is Gibbs dist. on

$G_Y$. with neighbour struc: $\eta_1 \sim \eta_2$ if $\begin{cases} \eta_1 \sim \eta_2 \text{ in } G_Y \\ \eta_1 \sim x \sim \eta_2. x \in G_X. \end{cases}$

ii) $X \mid Y$ is MRF on $G_X$. with neighbour struc. $N_X$.

Pf: i) $p(\eta) = \sum_{S_X} p(x, \eta)$. Written by def.

ii) $p(x \mid \eta) = p(x, \eta)/p(\eta)$.

(2) <u>Hidden Markov Chain</u> :

Consider :



$\theta = (S, A, B)$.  parameters :

i)  $S$ is initial dist of $X_0$.

ii)  $A_{ij} = P(X_{t+1}=j \mid X_t=i)$.  $B_{ij} = P(Y_t=j \mid X_t=i)$

$A = (A_{ij})_{u \times u}$.  $B = (B_{ij})_{u \times v}$.  prob. trans. Matrixs.

<u>Rmk</u>:  $u=1 \Rightarrow (Y_n)$ i.i.d.

$u=v$.  $B = I_u \Rightarrow (X_n)$ is Markov Chain

① <u>Likelihood</u> :

$L(\theta) = P_\theta(y_0 \cdots y_n)$.  density of observed data.

$L(\theta) = \sum_{x \in S_X} P_\theta(x, y)$.  $|S_X| = u^{n+1}$.

$\Rightarrow$ To calculate $L(\theta)$. We need sum up $u^{n+1}$ times.

<u>Denote</u>:  $\alpha_t(X_t) = P_\theta(X_t, y_0, y_1 \cdots y_t)$.  $\beta_t(X_t) = P_\theta(y_{t+1} \sim y_n \mid X_t)$

i) <u>Forward Prob.</u> :

$\alpha_0(X_0) = P_\theta(X_0, y_0) = S(X_0) B(X_0, y_0)$

$\alpha_{t+1}(X_{t+1}) = P_\theta(X_{t+1}, y_0 \cdots y_{t+1}) = \sum_{X_t} P_\theta(X_t, X_{t+1}, y_1 \cdots y_{t+1})$

$= \sum_{X_t \in S} \alpha_t(X_t) A(X_t, X_{t+1}) B(X_{t+1}, y_{t+1})$

$L_\theta(y) = \sum_{X_n \in S} \alpha_n(X_n)$.  obtained by iteratedly calculation

ii) <u>Backward Prob.:</u>

$$\beta_{t-1}(x_{b-1}) = P_\theta(\eta_{t}, \cdots \eta_n \mid x_{t-1}) = \sum_{S} P_\theta(x_t, \eta_t \sim \eta_n \mid x_{t-1})$$

$$= \sum_{x_t} P_\theta(x_{t-1}, x_t, \eta_t \sim \eta_n) / S(x_{t-1})$$

$$= \sum_{x_t} A(x_{t-1}, x_t) B(x_t, \eta_t) \beta_t(x_t)$$

$$L(\theta) = \beta_0(x_0) S(x_0) = P_\theta(\eta_1, \eta_2 \cdots \eta_n)$$

② <u>Maximize Likelihood:</u>

After calculating $L(\theta)$ given by $\theta \in \{S, A, B\}$.

We want to find $\hat{\theta}$ to max $L(\theta)$. Which

is best predictator for dist. of Hmm.

<u>Lemma.</u> For $p = (p_i)_i^k$, $z = (z_i)_i^k$ dist. on $(i)_i^k$.

We have: $\sum_i^k p_i \log p_i \geqslant \sum_i^k p_i \log z_i$.

<u>pf:</u> By $\sum p_i \log z_i / p_i \leqslant \sum p_i (z_i / p_i - 1) = 0$.

<u>Rmk:</u> Distance between dist. $p$. $z$:

$$D(p \| z) = \sum p_i \log p_i / z_i. \text{ is called}$$

kullback - Leibler distance.

To maximize $L_\theta = P_\theta(\eta) = \sum_x P_\theta(x, \eta) \iff$ Given

$Y = \eta$, maximize $P_\theta(x, \eta)$.

Next, we introduce EM Algorithm:

**prop.** If $\bar{E}_{\theta_0}(\log P_{\theta_1}(X,\eta)\mid\eta) > \bar{E}_{\theta_0}(\log P_{\theta_0}(X,\eta)\mid\eta)$.

Then: $P_{\theta_1}(\eta) > P_{\theta_0}(\eta)$.

**Pf:** $0 < E_{\theta_0}\left(\log\dfrac{P_{\theta_1}(X\mid\eta)}{P_{\theta_0}(X\mid\eta)}\mid Y=\eta\right)$

$= \sum_x P_{\theta_0}(x\mid\eta)\log P_{\theta_1}(\eta)/P_{\theta_0}(\eta) - \sum_x P_{\theta_0}(x\mid\eta)\log\dfrac{P_{\theta_1}(x\mid\eta)}{P_{\theta_0}(x\mid\eta)}$

$= \log P_{\theta_1}(\eta)/P_{\theta_0}(\eta) - \sum\square \le \log P_{\theta_1}(\eta)/P_{\theta_0}(\eta)$

Note that: $P_\theta(x,\eta) = \mathcal{S}(x_0)\prod_0^{n-1} A(x_t,x_{t+1})\prod_0^n B(x_t,\eta_t)$

$\Rightarrow \log P_\theta(x,\eta) = \log \mathcal{S}(x_0) + \sum\log A(x_t,x_{t+1}) + \sum\log B(x_t,\eta_t)$

1') Randomly choose $\theta_0(\mathcal{S}_0,A_0,B_0)$

2') Choose $\theta_1=(\mathcal{S}_1,A_1,B_1)$ maximizes $f(\theta) = \bar{E}_{\theta_0}(\log P_\theta(x,\eta)\mid\eta)$

$= \sum_i P_{\theta_0}(X_0=i\mid\eta)\log\mathcal{S}(i) + \sum_{t=0}^{n-1}\sum_{i,j} P_{\theta_0}(X_t=i,X_{t+1}=j\mid\eta)\log A(i,j)$

$+ \sum_{t=0}^n \sum_i P_{\theta_0}(X_t=i\mid\eta)\log B(i,\eta_t) \stackrel{\triangle}{=} A_1 + A_2 + A_3.$

---

For $A_1$: Choose $\mathcal{S}_1(i) = P_{\theta_0}(X_0=i\mid\eta)$, by Lemma.

the prob. condition on current data.

For $A_2$: $\sum_i\sum_j\left(\sum_t P_{\theta_0}(X_t=i,X_{t+1}=j\mid\eta)\right)\log A(i,j))$

choose $A_1(i,j) = \left(\sum_t^{n-1} P_{\theta_0}(X_t=i,X_{t+1}=j\mid\eta)\right)\Big/\sum_j\left(\sum_0^{n-1}\square\right)$

$(\hat{A}(i,j) = \dfrac{\sum_t I(X_t=i,X_{t+1}=j)}{\sum_t I(X_t=i)} \approx p(X_t=i\mid X_{t+1}=j))$

For $A_3$: By Lemma analogously, Choose:

$B_1(i,j) = \sum_{t:\eta_t=j} P_{\theta_0}(X_t=i\mid\eta)\Big/\sum_j\sum_{t:\eta_t=i} P_{\theta_0}(X_t=i\mid\eta)$

$(\hat{B}(i,j) = \dfrac{\sum_t^n I(X_t=i,Y_t=j)}{\sum_t^n I(X_t=i)} \approx p(Y_t=i\mid X_t=j))$

3') Calculate $\theta_1 = (S_1, A_1, B_1)$

Consider $Y_t(i,j) = P_{\theta_0}(X_t = i, X_{t+1} = j \mid \eta)$ which

can express $\theta_1$. $\Longleftrightarrow P_{\theta_0}(X_t, X_{t+1}, \eta)$. since $P_{\theta_0}(\eta)$

can be calculated by forward prob.

$$P_{\theta_0}(X_t, X_{t+1}, \eta) = P_{\theta_0}(\eta_0^t, X_t, X_{t+1}, \eta_{t+1}^n)$$

$$= \alpha_t(X_t) A_0(X_t, X_{t+1}) P_{\theta_0}(\eta_{t+1}^n \mid X_{t+1})$$

$$P_{\theta_0}(\eta_{t+1}^n \mid X_{t+1}) = P_{\theta_0}(\eta_{t+1}, \eta_{t+2}^n \mid X_{t+1})$$

$$= B_0(X_{t+1}, \eta_{t+1}) \beta_{t+1}(X_{t+1})$$

$$\Rightarrow Y_t(i,j) = \alpha_t(X_t) A_0(X_t, X_{t+1}) B_0(X_{t+1}, \eta_{t+1}) \beta_{t+1}(X_t) \Big/ \square$$

$$\square = \sum_{i,j} \alpha_t(i) A_0(i,j) B(j, \eta_{t+1}) \beta_{t+1}(j).$$

Def: $S_1(i) = \sum_j Y_0(i,j)$

$$A_1(i,j) = \sum_{t=0}^{n-1} Y_t(i,j) \Big/ \sum_{l} \sum_{t=0}^{n-1} Y_t(i,l)$$

$$B_1(i,j) = \frac{\sum_{t=0}^{n-1} \sum_{l} Y_t(i,l) I_{\{\eta_t = j\}} + \sum_{m} Y_n(m,i) I_{\{\eta_n = j\}}}{\sum_{t=0}^{n-1} \sum_{l} Y_t(i,l) + \sum_{m} Y_n(m,i)}$$

4') Replace $\theta_0$ by $\theta_1$. Repeat the procedure.