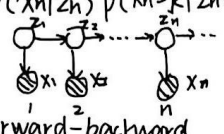


**HMM**  $\theta = (\pi, A, B)$   
 Latent  $Z_n$  (时刻  $n$ )  
 $A = [a_{ij}]$  由状态  $i$  转移到  $j$  的概率  
 $B = [b_{jk}]$  发射概率, 观测变量的概率  
 $p(x_n | z_n) p(z_n = k | z_{n-1}) = b_{jk}$   
  
 $z_n(i)$  时, 模型处于  $i$   
 且观测序列  $x_1, x_2, \dots, x_n$   
 已经发生时条件下  
 状态序列的概率

**forward-backward**  
 前  $d_n(i) = P(x_1, \dots, x_n, z_n = i | \theta)$   
 $d_1(i) = \pi_i \cdot b_i(x_1)$  sum-product  
 $d_{n+1}(j) = (\sum_{i=1}^N d_n(i) a_{ij}) b_j(x_{n+1})$   
 后  $B_n(i) = P(x_{n+1}, x_{n+2}, \dots, x_N | z_n = i, \theta)$   
 $B_N(i) = 1$  (一定出现) max-product  
 $B_{n+1}(i) = \sum_{j=1}^N a_{ij} \cdot b_j(x_{n+1}) \cdot B_n(j)$   
 $p(z_n = i | x_1, \theta) = \frac{d_n(i) B_n(i)}{P(x | \theta)} \rightarrow \sum_{i=1}^N d_n(i) B_n(i)$   
 $P(z_{n+1}, z_n | x_1, \theta) = \frac{d_{n+1}(j) B_j(x_n) a_{ji} B_n(i)}{P(x | \theta)}$

**MDP**  $\langle S, R, P, \gamma \rangle$   
 state  $x$ , next  $x'$ , action  $u$  policy  $\pi(u|x)$   
 转移概率  $p(x'|u,x)$  (与  $t$  无关)  
 reward  $r(x,u)$  discount  $\gamma \in (0,1)$   
 $R_t^\pi(x_t) = E[\sum_{\tau=0}^{\infty} \gamma^\tau r_{t+\tau} | u_{t+\tau} = \pi(x_{t+\tau})]$   
 指标  $\pi^* = \arg \max_{\pi} R_t^\pi(x_t)$   
 additive  $R(x_t, x_{t+1}, \dots) = r(x_t) + \gamma R(x_{t+1}) + \gamma^2 R(x_{t+2}) + \dots$   
 state value  $V(x)$  从  $x$  状态开始  $R$  的期望  
 $V(x) = E_\pi[\sum_{\tau=0}^{\infty} \gamma^\tau R_{t+\tau} | x_t = x] = r(x,u) + \gamma \sum_{x'} p(x'|x,u) V(x')$   
 $Q(s,u) = E_\pi[r_{t+1} + \gamma V(x_{t+1}, u_{t+1}) | x_t = s, u_t = u]$   
 $V^*(x|x) = \sum_{u \in A} \pi(u|x) Q(x,u)$   
 $Q^*(x) = \sum_{u \in A} \pi(u|x) r(x,u)$   
 $V(s,u) = r(s,u) + \gamma \sum_{x'} p(x'|s,u) \cdot V(x')$   
 $Q(s) = \sum_{u \in A} (u|s) q_\pi(s,u) \rightarrow$   
 $V(s) = \sum_{u \in A} \pi(u|s) (r(s,u) + \gamma \sum_{x'} p(x'|s,u) \cdot V(x'))$   
 $\rightarrow \sum_{x'} \gamma V(x') = 0 \quad V(s) = \max_u V(s,u)$   
 $V^*(s,a) = r(s,a) + \gamma \sum_{x'} p(x'|s,a) \max_{a'} V(x',a')$

$u$ : regression classification  
 detection  
 KNN, SVM, DT, NN, DNN  
 $u$ : Clustering, data Dimension reduction  
 K-means, GMM, PCA, ICA, NMF, GAN

**K-Means**  $M_k$   
 1. 随机找  $K$  点, 2. 把每个点分给最近的中心  
 3. 重新计算中心点  
 终止  $J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - M_k\|^2$   
 I.  $\min J \rightarrow 0$  or II.  $\|r_{nk}\| \min M_k$   
 $M_k = \frac{\sum_{n \in k} r_{nk} x_n}{\sum_{n \in k} r_{nk}}$   $\rightarrow$  混合系数

**GMM**  $p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$   
 对每个观测数据  $x_n$  有 latent  $z_n$  聚类标签  
 $E: \gamma(z_n) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)} = p(z_n = k | x_n)$   
 给定  $x_n$  属于第  $k$  个高斯分布的概率  
 $p(x | \pi, \mu, \Sigma) = \prod_{n=1}^N \left[ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right]$   
 maximum 指导

$M_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_n) x_n \quad N_k = \sum_{n=1}^N \gamma(z_n)$   
 $\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_n) (x_n - M_k^{new})(x_n - M_k^{new})^T$   
 $\pi_k^{new} = \frac{N_k}{N}$   
 complete data  $\{x_i, z_i\}$   
 $p(x, z | \mu, \pi) = \prod_{n=1}^N \prod_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)^{z_{nk}}$   
 $E(z_{nk}) = \gamma(z_n)$   
 $E_z[\ln p(x, z | \mu, \pi)] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_n) \{ \ln \pi_k + \ln N(x_n | \mu_k, \Sigma_k) \}$

**Mixture of Bernoulli**  
 $p(x|z, \mu) = \prod_{k=1}^K p(x_k | \mu_k)^{z_k}, p(z | \pi) = \prod_{k=1}^K \pi_k^{z_k}$   
 $\ln p(x | \mu, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k p(x_n | \mu_k) \right\}$   
 $\ln p(x | z | \mu, \pi) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \sum_{j=1}^D [x_{nj} \ln \mu_{kj} + (1-x_{nj}) \ln (1-\mu_{kj})] \}$   
 $E_z: \gamma(z_{nk}) = E(z_{nk}) = \frac{\sum_{j=1}^D z_{nj} \mu_{kj} p(x_n | \mu_j)}{\sum_{j=1}^D [\pi_j p(x_n | \mu_j)]^{z_{nj}}}$   
 $= \frac{\pi_k p(x_n | \mu_k)}{\sum_{j=1}^D \pi_j p(x_n | \mu_j)}$   
 $E_z[\ln p(x, z | \mu, \pi)] = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \dots \}$   
 $M: N_k = \sum_{n=1}^N \gamma(z_{nk}), \pi_k = \frac{N_k}{N}$   
 $\bar{x}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) x_n \quad \mu_k = \bar{x}_k$   
 (E) 隐变量后验  $\gamma(z)$   
 (M) 更新模型参数 (根据  $E(z)$ )

**EM通用算法** 求  $p \ln p(x|\theta)$   
 $\ln p(x|\theta) = L(q, \theta) + KL(q||p)$  (非负)  
 $L(q, \theta) = \sum_{z=1}^Z q(z) \ln \left( \frac{p(x, z | \theta)}{q(z)} \right)$   
 $KL(q||p) = - \sum_{z=1}^Z q(z) \ln \left( \frac{q(z)}{p(x, z | \theta)} \right) \geq 0$   
 两概率分布差异  
 $\Rightarrow$  最大化  $L(q, \theta)$   $\rightarrow$  对 latent 变量后验评估下界  
 $E: \text{Fix } \theta^{old} \downarrow L(q, \theta^{old}) \max_{\theta} KL(q||p) = 0$   
 $M: \text{Fix } q(z) \downarrow L(q, \theta^{old}) \max_{\theta} KL(q||p) = 0$   
 $L(q, \theta) = \sum_{z=1}^Z q(z) \ln p(x, z | \theta) - \sum_{z=1}^Z q(z) \ln q(z)$   
 其  $q(z) = p(z | x, \theta^{old})$   
 $\Rightarrow L(q, \theta) = Q(\theta, \theta^{old}) + \text{const}$

**EM for Linear Regression**  
 $\ln p(w, b, \beta) = \ln p(w, b) + \ln p(\beta)$   
 $E \Rightarrow E[\ln p(w, b, \beta)] = \sum_{n=1}^N \ln \left( \frac{1}{2\pi} \right) - \frac{1}{2} E[\ln p(w, b, \beta)]$   
 $+ \sum_{n=1}^N \ln \left( \frac{1}{2\pi} \right) - \frac{1}{2} \sum_{n=1}^N E[\ln p(w, b, \beta)]$   
 $M \rightarrow d = \frac{M}{M} = \frac{M}{M}$   
 $\frac{1}{\beta} = \frac{1}{N} \sum_{n=1}^N (y_n - M^T \phi_n)^2$   
 补 MDP 求解  
 Bellman  $V^*(s) = \max_a \sum_{s'} p(s'|s,a) [R(s,a,s') + \gamma V^*(s')]$   
 ① value iteration  
 $V^*(s) = r(s,a) + \gamma V^*(s') p(s'|s,a)$  选  $a$  向最大的  
 ② policy iteration  
 $V^* \rightarrow \pi^*$

**HMM + EM**  
 $E\text{-step } d_t(z) = \sum_{z'} d_{t-1}(z') A_{zz'} p(x_t | z') \beta_{t-1}(z')$   
 $\beta_t(z) = \sum_{z'} A_{zz'} p(x_{t+1} | z') \beta_{t+1}(z')$   
 状态边缘概率  $z$   
 $\gamma_t(z) = p(z_t = z | x_1, \theta) \propto d_t(z) \beta_t(z)$   
 状态转移边缘概率  
 $\xi_t(z, z') = p(z_t = z, z_{t+1} = z' | x_1, \theta) \propto d_t(z) A_{zz'} p(x_{t+1} | z') \beta_{t+1}(z')$   
 $M\text{-step} =$  更新模型参数  $\theta = \{\pi, A, B\}$

$\pi(z) = \gamma_1(z)$   
 $A_{zz'} = \frac{\sum_{t=1}^T \xi_t(z, z')}{\sum_{t=1}^T \gamma_t(z)}$   
 $B_k = \frac{\sum_{t=1}^T \sum_{n=1}^N \gamma_t(z_n) x_{nk}}{\sum_{t=1}^T \gamma_t(z_n)}$   
 Ex.  $p(x|\theta) = \theta^x (1-\theta)^{1-x} \rightarrow \theta^x$   
 $p(x, z | \theta) = \prod_{i=1}^I [p(z_i) p(x_i | \theta)]^{z_i}$   
 $M: Q(\theta | \theta^{old}) = E_{z|x, \theta^{old}} [\log p(x, z | \theta)]$   
 对  $\theta$  指导  
**MAP + EM**  
 $E\text{-step}$   
 $M: Q(\theta | \theta^{old}) = E[\log p(x, z | \theta)] + \ln p(\theta)$

**SVM**  $y(x) = w^T \phi(x) + b$   
 $\ln y(x_n) > 1$   
 $\text{dis}(\text{point} \rightarrow \text{line}) = \frac{\ln y(x_n)}{\|w\|}$   
 目标:  $\min \|w\|^2$  s.t.  $\ln(w^T \phi + b) \geq 1$   
 Lagrange  
 $L(w, b, \alpha) = \frac{\|w\|^2}{2} - \sum_{n=1}^N \alpha_n (\ln(w^T \phi(x_n) + b) - 1)$   
 $\Rightarrow w = \sum_{n=1}^N \alpha_n \phi(x_n)$   
 $\sum_{n=1}^N \alpha_n \ln \phi(x_n) = 0$   
 Dual  $\tilde{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m \ln \phi(x_n) \ln \phi(x_m)$   
 s.t.  $\alpha_n \geq 0 \quad \sum_{n=1}^N \alpha_n \ln \phi(x_n) = 0$   
 quadratic programming.  
 solve an  
 Max Margin Classifier  $\lambda > 0$   
 $\arg \min \sum_{n=1}^N E_{\phi}(\ln y(x_n) - 1) + \lambda \|w\|^2$   
 $\rightarrow \begin{cases} 0 & \text{if } z \geq 0 \\ \infty & \text{else} \end{cases}$   
 Soft / overlap  
 $\arg \min \sum_{n=1}^N \ln y(x_n) \geq 1 - \epsilon_n$   
 $\min \frac{\|w\|^2}{2} + C \sum_{n=1}^N \epsilon_n$   
 $L(w, b, \alpha) = \frac{\|w\|^2}{2} + C \sum_{n=1}^N \epsilon_n - \sum_{n=1}^N \alpha_n (\ln y(x_n) - 1 + \epsilon_n)$

$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum \alpha_n \ln \phi(x_n)$   
 $\frac{\partial L}{\partial b} = 0 \Rightarrow \sum \alpha_n \ln = 0$   
 $\frac{\partial L}{\partial \epsilon_n} = 0 \Rightarrow \alpha_n = C - \mu_n$   
**KKT Condition**  
 $\mu_n \geq 0 \quad \epsilon_n \geq 0 \quad \mu_n \epsilon_n = 0$   
 $\alpha_n \geq 0 \quad \ln y(x_n) \geq 1 - \epsilon_n \quad \alpha_n (\ln y(x_n) - 1 + \epsilon_n) = 0$   
 $\frac{\partial L}{\partial w} = \frac{\partial L}{\partial b} = \frac{\partial L}{\partial \epsilon_n} = 0$



# NN

error funct: cross entropy

$$E(w) = - \sum_{n=1}^N (t_n \ln y_n + (1-t_n) \ln(1-y_n))$$

y = sigmoid

regression y = a

$$E(w) = \frac{1}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2$$

multiclass

$$E(w) = - \sum_n \sum_k t_{nk} \ln y_{k|n}(x_n, w)$$

$$w_{ji}^{t+1} = w_{ji}^t + \Delta w_{ji}^t = w_{ji}^t - \eta \nabla E(w_{ji}^t)$$

$$E = \frac{1}{2} \sum_k (y_k - t_k)^2$$

$$a_j = \sum_i w_{ji} z_i \quad z_j = h(a_j) \quad \frac{\partial a_j}{\partial w_{ji}} = z_i$$

$$a_k = \sum_j w_{kj} z_j \quad y_k = \sigma(a_k)$$

$$\frac{\partial E}{\partial a_j} = \sum_k \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial a_j} = \sum_k (y_k - t_k) w_{kj} \cdot h'(a_k)$$

$$E_n = - \sum_k t_{nk} \ln y_{k|n} + (1-t_{nk}) \ln(1-y_{k|n})$$

$$\frac{\partial E_n}{\partial y_k} = \frac{y_k - t_{nk}}{y_k(1-y_k)} \frac{\partial y_k}{\partial a_k} y_k(1-y_k), y_k = \sigma(a_k)$$

$$\frac{\partial E_n}{\partial w_{kj}} = \delta_k z_j \quad \frac{\partial E_n}{\partial w_{ji}} = \delta_j z_i$$

$$H = \nabla^2 E = \sum_{n=1}^N \nabla y_n \nabla y_n + \sum_{n=1}^N (y_n - t_n) \nabla \nabla y_n$$

## NN-regression

$$E(w) = - \ln p(w|t) = \frac{1}{2} w^T w + \frac{1}{2} \sum_{n=1}^N (y(x_n, w) - t_n)^2 + c$$

$$\nabla E(w) = \Delta w + \beta \sum_{n=1}^N (y_n - t_n) \nabla w y(x_n, w)$$

$$A = \nabla \nabla E(w) = \Delta I + \beta H$$

$$w_{map} \leftarrow w_{new} = w_{old} - A^{-1} \nabla E(w)$$

$$q(w) = N(w|w_{map}, A^{-1})$$

## NN-classification

$$E(w) = - \ln p(w|t) = \frac{1}{2} w^T w - \sum_{n=1}^N (t_n \ln y_n + (1-t_n) \ln(1-y_n))$$

$$\nabla E(w) = \Delta w + \sum_{n=1}^N (y_n - t_n) g_n$$

$$A = \nabla \nabla E(w) = \Delta I + H$$

$$w_{map} \leftarrow w_{new} = w_{old} - A^{-1} \nabla E(w)$$

$$p(t|x, D) = \int p(t|x, w) q(w|D) dw$$

Evaluation

$$p(D) = \int p(D|\theta) p(\theta) d\theta$$

$$\ln p(D) = \ln p(D|\theta_{map}) + \ln p(\theta_{map}) + \sum_{n=1}^M \ln p(t_n)$$

$$\text{where } A = - \nabla \nabla \ln p(D|\theta_{map}) p(\theta_{map}) - \frac{1}{2} M I$$

Discrimination

$$y(x) = w^T x \quad E(w) = \frac{1}{2} (XW - T)^T (XW - T)$$

$$\tilde{w} = (X^T X)^{-1} X^T T$$

Fisher

$$y = w^T x \quad w \propto m_2 - m_1$$

$$S_1^2 = \sum_{n \in C_1} (y_n - m_1)^2 \quad S_2^2 = \sum_{n \in C_2} (y_n - m_2)^2$$

$$J(w) = \frac{(m_2 - m_1)^2}{S_1^2 + S_2^2} \quad \text{objective} = \frac{w^T S_2 w}{w^T S_1 w}$$

$$S_B = (m_2 - m_1)(m_2 - m_1)^T$$

$$S_W = \sum_{n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{n \in C_2} (x_n - m_2)(x_n - m_2)^T$$

$$w \propto S_W^{-1} (m_2 - m_1)$$

# Reinforcement Learning for MDP

1. What is the Bellman equation? How to solve the Bellman equation?

$v_\pi(s) = E[r_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$ . Analytical sol:  $V = (I - \gamma P)^{-1} R$  using Gaussian, can also be solved by dynamic programming, Monte-Carlo method, and temporal difference method.

What are the differences between policy iteration and value iteration? What are their advantages and disadvantages respectively?

• **Policy iteration:** Policy evaluation + policy enhancement. (Attach equations to explain) Policy evaluation in the policy iteration uses the Bellman expectation equation to obtain the state value function of a policy, which is a dynamic programming process; while the value iteration directly uses the Bellman optimization equation for dynamic programming to obtain the final optimal state value.

• **Value iteration:** (Attach equations to explain) Only one round of value updates in the policy evaluation, then do policy enhancements based directly on the updated values. Finally use  $\epsilon$ -policy (refer to MDP) to obtain the optimal policy.

Policy Iteration	Value Iteration
Complex algorithm	Simple algorithm
Cheaper to compute	Expensive to compute
Faster convergence	Slower convergence

3. What is the model-free reinforcement learning? How to achieve the model-free reinforcement learning? Please use specific examples to illustrate your points.

Methods that agent can only learn through data that comes from interacting with the environment is called model-free reinforcement learning (Sarsa and Q-learning).

Application: Complex environment like E-sport games.

4. What are the differences between on-line and off-line RL? What are their advantages and disadvantages respectively? Please use specific examples to illustrate your points.

**Online RL:** Learning from the data that sampled from the environment by the current policy, and once the policy is updated, the data points will be deprecated. (E.g. Sarsa, policy is greedy)

**Offline RL:** Use experience recall pool to store and reuse the historical samples (E.g. Q-learning). Offline RL can better make use of historical data and obtain lower complexity of sample, i.e. use smaller amount of data to achieve convergence.

$$\phi = \begin{pmatrix} \phi_0(x) \\ \phi_1(x) \end{pmatrix}$$

$$tr(ABC) = tr(CAB)$$

$$= tr(BCA)$$

$$\frac{\partial tr(AB)}{\partial A} = B^T$$

$$\frac{\partial tr(AX^T S^{-1})}{\partial S^{-1}} = X X^T$$

$$\frac{\partial \ln |A|}{\partial A} = (A^T)^{-1}$$

$$\frac{\partial \ln |\Sigma|}{\partial \Sigma^{-1}} = \Sigma$$

$$y = Ax + v \quad x \sim N(0, \Sigma) \quad v \sim N(0, \Sigma_v)$$

$$x = m + u \quad u \sim N(0, I) \quad p(x|y) = N(x|m, I)$$

$$\Rightarrow L^T = A^T Q^T A + \Sigma^{-1}$$

$$L m = A^T Q^T y + \Sigma^{-1} \mu$$

$$p(y|x) = N(y|Ax, \Sigma)$$

$$p(y) = \int p(y|x) p(x) dx = N(y|A \mu, A \Sigma A^T + \Sigma_v)$$

$$p(y) = \int p(y|x) p(x) dx = N(y|A \mu, A \Sigma A^T + \Sigma_v)$$

$$t = y(x, w) + \epsilon \rightarrow (0, \beta)$$

$$p(t|x, w, \beta) = N(t|y(x, w), \beta^{-1})$$

$$\text{likelihood} = \prod_{n=1}^N N(t_n | y(x_n, w), \beta^{-1})$$

$$E_D(w) = \frac{1}{2} \sum (t_n - w^T \phi(x_n))^2 \quad w_{ML} = (\Phi^T \Phi)^{-1} \Phi^T T$$

$$E_D(w) + \lambda E_M(w) \quad w = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T T$$

Bayesian Linear regression

$$p(w) = N(w|m_0, S_0)$$

$$p(w|t) = N(w|m_N, S_N)$$

$$S_N^{-1} m_N = \beta \Phi^T t + S_0^{-1} m_0$$

$$S_N^{-1} = \beta \Phi^T \Phi + S_0^{-1}$$

$$\text{learn } p(\theta | \text{train, model}) = p(t|m, \theta) p(\theta|m)$$

$$\text{prediction } p(t | \text{test, train, model}) = \int p(t | \text{test, } m, \theta) p(\theta | \text{train, } m) d\theta$$

$$\text{evaluation } p(t | \text{train, } m) = \int p(t | \text{train, } m, \theta) p(\theta | m) d\theta$$

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu))$$

$$\text{cross-entropy} \quad \text{likelihood} \quad \text{ML与MAP 考虑先验} \quad \text{EM局部 MDP stationary 时间独立}$$