

Example: Polynomial Regression

The polynomial regression model, jointly showing also a new input value \tilde{x} together with the corresponding model prediction \hat{t} .

$$p(\tilde{t} | \mathbf{x}_{1:N}, \mathbf{w}, \alpha, \sigma^2) = \left[\prod_{n=1}^N p(x_n | \mathbf{x}_n, \mathbf{w}, \alpha^2) \right] \times p(\mathbf{w} | \mathbf{x}) p(\tilde{x} | \mathbf{x}, \mathbf{w}, \sigma^2)$$

参数
Plate (replication)
Observed Variable
Latent Variable
hyperparameter
可能为零
D-separation
随机变量
确定变量
可观察变量
随机变量
参数
重复出现次数
noise (sign of)
D-separation
All paths from any node in A (node set) to any node in B (node set) are disconnected given $C \not\leftrightarrow A$ is said to be D-separated from B by C ($A \perp\!\!\!\perp B | C$)

Generating process of GMM:

- First generate cluster index: $\mathbf{z} \sim (\pi_1, \dots, \pi_k)$
- Then generate data: $\mathbf{x} \sim \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$.

Mixture Distribution

A probability density $p(\mathbf{x})$ represents a mixture distribution

- where $w_i \geq 0, \sum_{i=1}^k w_i = 1$, and each p_i is a probability density.

$$p(\mathbf{x}) = \sum_{i=1}^k w_i p_i(\mathbf{x})$$

More constructively, let S be a set of probability distributions:

- Choose a distribution randomly from S .
- Sample \mathbf{x} from the chosen distribution.
- Then \mathbf{x} has a mixture distribution.

Expectation-Maximization (EM): Algorithm

Choose initial θ^{old}

Expectation Step

Let $q^*(\mathbf{z}) = p(\mathbf{z} | \mathbf{x}, \theta^{\text{old}})$. [q^* gives best lower bound at θ^{old}]

$$J(\theta) := L(q^*, \theta) = \sum_{\mathbf{z}} q^*(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q^*(\mathbf{z})}$$

Maximization Step

$\theta^{\text{new}} = \operatorname{argmax}_{\theta} J(\theta)$ You can use SGD

Go to the Expectation Step, until converged.

EM for MAP

• Suppose we have a prior $p(\theta)$.

We still model $p(Y = y)$ as Bernoulli distribution: $Bernoulli(\phi)$.

Now we use MLE to find the best parameter estimation:

$$\ell(\phi, \mu_+, \mu_-) = \log \prod_{i=1}^n p(x_i | y_i; \mu_+, \mu_-, \Sigma)$$

$$= \log \prod_{i=1}^n p(x_i | y_i; \mu_+, \mu_-, \Sigma) + \log \prod_{i=1}^n p(y_i | \phi)$$

The computing process is very similar to the process of Gaussian.

- The main difference is that μ_+, μ_- are different.

$$\phi = \frac{\sum_{i=1}^n \mathbf{1}(y_i=+1)}{n}, \mu_+ = \frac{\sum_{i=1}^n \mathbf{1}(y_i=+1)x_i}{\sum_{i=1}^n \mathbf{1}(y_i=+1)}, \mu_- = \frac{\sum_{i=1}^n \mathbf{1}(y_i=-1)x_i}{\sum_{i=1}^n \mathbf{1}(y_i=-1)}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_+)(x_i - \mu_-)^T$$

Gaussian Mixture Model

Parameters of

Gaussian Mixture Model (GMM):

- Cluster probabilities: $\pi = (\pi_1, \dots, \pi_k)$.

- Cluster means: $\mu = (\mu_1, \dots, \mu_k)$.

- Cluster covariance matrices: $\Sigma = (\Sigma_1, \dots, \Sigma_k)$

$$\text{Density: } p(\mathbf{x}) = \sum_{z=1}^k \pi_z \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$$

GMM: M-Step

$$\begin{aligned} \operatorname{argmax}_{\theta} \mathcal{L}(q^*, \theta) &= \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^*(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q^*(\mathbf{z})} \\ \theta = \{\pi, \boldsymbol{\mu}, \boldsymbol{\Sigma}\} &= \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} p(\mathbf{z} | \mathbf{x}, \theta^{\text{old}}) \log(p(\mathbf{x}, \mathbf{z} | \theta)) \end{aligned}$$

• So we have the loss function for Gaussian Mixture Model parameters:

$$\text{By MLE } \operatorname{argmax}_{\pi, \boldsymbol{\mu}, \boldsymbol{\Sigma}} \sum_{i=1}^n \sum_{j=1}^k \gamma_i^j \log [\pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)]$$

• Let $n_c = \sum_{i=1}^n \gamma_i^c$ be the number of points soft-assigned to cluster c .

$$\pi_c^{\text{new}} = \frac{n_c}{n}, \boldsymbol{\mu}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^n \gamma_i^c \mathbf{x}_i, \boldsymbol{\Sigma}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^n \gamma_i^c (\mathbf{x}_i - \boldsymbol{\mu}_c^{\text{new}})(\mathbf{x}_i - \boldsymbol{\mu}_c^{\text{new}})^T$$

EM for GMM: Overview

Initialize parameters $\pi, \boldsymbol{\Sigma}, \boldsymbol{\mu}$.

$$\gamma_i^j = \frac{\pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{c=1}^k \pi_c \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)}$$

E-step. Evaluate all responsibilities using current parameters:

M-step. Re-estimate the parameters using the responsibilities:

$$\pi_c^{\text{new}} = \frac{n_c}{n}, \boldsymbol{\mu}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^n \gamma_i^c \mathbf{x}_i, \boldsymbol{\Sigma}_c^{\text{new}} = \frac{1}{n_c} \sum_{i=1}^n \gamma_i^c (\mathbf{x}_i - \boldsymbol{\mu}_c^{\text{new}})(\mathbf{x}_i - \boldsymbol{\mu}_c^{\text{new}})^T$$

Repeat E-step and M-step until log-likelihood converges.

EM and Variational Methods

• When E-step is difficult:

- Hard to take expectation w.r.t. $q^*(\mathbf{z}) = p(\mathbf{z} | \mathbf{x}, \theta^{\text{old}})$

- For example, hierarchical latent variable models (next lectures).

• Solution: Restrict to distributions \mathcal{Q} that are easy to work with.

• The evidence lower bound (ELBO) now looser:

$$q^* = \operatorname{argmin}_{\mathcal{Q}} \text{KL}[q(\mathbf{z}), p(\mathbf{z} | \mathbf{x}, \theta^{\text{old}})]$$

• Find an easy-to-work variational distribution q^* to approximate the inference distribution $p(\mathbf{z} | \mathbf{x}, \theta^{\text{old}})$.

- This group of methods are called variational methods.

Dirichlet-Multinomial Mixture Model

• Assume that we have a vocabulary of V words.

• Generative Process

$$\begin{aligned} \theta &\sim \text{Dir}(\alpha) && \text{draw distribution over words} \\ \text{For each word } n \in \{1, \dots, N\} && \text{draw word} \\ w_n &\sim \text{Mult}(1, \theta) && \text{draw word} \end{aligned}$$

Categorical Distribution

• Still can use our evidence lower bound on $\log p(\mathbf{x}, \theta)$.

$$J(\theta) := \mathcal{L}(q^*, \theta) = \sum_{\mathbf{z}} q^*(\mathbf{z}) \log \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q^*(\mathbf{z})}$$

• Maximization step becomes $\theta^{\text{new}} = \operatorname{argmax}_{\theta} [J(\theta) + \log p(\theta)]$

GMM: E-Step

• Denote probability (responsibility) that \mathbf{x}_i comes from cluster j by

$$\gamma_i^j = p(\mathbf{z} = j | \mathbf{x} = \mathbf{x}_i, \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

- The vector $(\gamma_1^1, \dots, \gamma_N^k)$ is exactly the soft assignment for \mathbf{x}_i .

• From probabilistic computation:

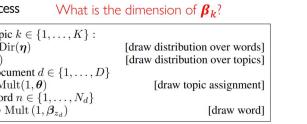
$$\gamma_i^j = p(\mathbf{z} = j | \mathbf{x}_i) = \frac{p(\mathbf{z} = j, \mathbf{x}_i)}{p(\mathbf{x}_i)} = \frac{\pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{c=1}^k \pi_c \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)}$$

• If we know $\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j, \pi_j$ for all clusters $j = 1, \dots, k$, then easy to compute:

$$\gamma_i^j = \frac{\pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{c=1}^k \pi_c \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)}$$

Dirichlet-Multinomial Mixture Model

What is the dimension of $\boldsymbol{\theta}$? What is the dimension of $\boldsymbol{\beta}_k$?



• Example corpus

the he is
w11 w12 w13

the and the
w21 w22 w23

the she she is
w31 w32 w33 w34

Document 1 Document 2 Document 3

Topics
Documents
Topics
Documents

Mixture
Admixture

There should be multiple topics in a document:
Each document can be modeled by a new distribution of topics!

Latent Dirichlet Allocation (LDA)

Each document has a single topic \mathbf{z}

$$\begin{aligned} \text{Dirichlet-Multinomial} &\text{Mixture Model} \\ \text{For each topic } k \in \{1, \dots, K\}: & \text{draw distribution over words} \\ \beta_k \sim \text{Dir}(\eta) & \text{draw distribution over topics} \\ \text{For each document } d \in \{1, \dots, D\}: & \text{draw topic assignment} \\ z_d \sim \text{Mult}(1, \theta) & \text{draw word} \\ \text{For each word } n \in \{1, \dots, N_d\}: & \text{draw word} \end{aligned}$$

• We solve one of the terms as an example:

$$E_{\theta}[\log p(\mathbf{z} | \boldsymbol{\alpha})]$$

$$= E_{\theta}[\log \exp \left(\sum_{i=1}^K (\alpha_i - 1) \log \theta_i \right) + \log \Gamma(\sum_{i=1}^K \alpha_i) - \sum_{i=1}^K \log \Gamma(\alpha_i)]$$

$$= \left(\sum_{i=1}^K (\alpha_i - 1) E_{\theta}[\log \theta_i] \right) + \log \Gamma(\sum_{i=1}^K \alpha_i) - \sum_{i=1}^K \log \Gamma(\alpha_i)$$

$$= \left(\sum_{i=1}^K (\alpha_i - 1) (\Psi(\gamma_i) - \Psi(\sum_{j \neq i} \gamma_j)) \right) + \log \Gamma(\sum_{i=1}^K \alpha_i) - \sum_{i=1}^K \log \Gamma(\alpha_i)$$

• Ψ is the first derivative of the log Gamma function Γ .

• Iterate the following steps (Descend each variable iteratively):

(1) For each topic k and term n :

$$\gamma_{k,n}^{(t+1)} = \eta + \sum_{d=1}^D \sum_{n=1}^N \mathbb{1}(w_{d,n} = v) \phi_{d,k}^{(t)}$$

(2) For each document d :

(a) Update $\phi_{d,k}^{(t)}$:

$$\phi_{d,k}^{(t+1)} = a_k + \sum_{n=1}^N \phi_{d,n,k}^{(t)}$$

(b) For each word n , update $\phi_{d,n}^{(t)}$:

$$\phi_{d,n}^{(t+1)} \propto \exp \left(\psi_{d,n}^{(t+1)} + \psi_{d,n}^{(t+1)} - \psi \left(\sum_{k=1}^K \phi_{d,k}^{(t+1)} \right) \right),$$

where ψ is the digamma function, the first derivative of the log gamma function.

$$\sum_{d=1}^D \phi_{d,n}^{(t+1)} = 1$$

• Markov chain is defined via a state transition matrix:

- Assume the state is discrete with K possible outcomes:

$$A_{ij} = P(x_{t+1,j} = 1 | x_{t,i} = 1), \quad i, j \in \{1, 2, \dots, K\}$$

State Transition Matrix

• If we unfold the state nodes → states edges → state transition matrix

If we unfold the state transition diagram over time, we obtain a lattice.

In HMMs, Viterbi algorithm is used for minimizing sequence error rate.

利用贝叶斯决策准则把误差最小化转化为最大后验概率

Suppose E-step is difficult:

- Hard to take expectation w.r.t. $q^*(\mathbf{z}) = p(\mathbf{z} | \mathbf{x}, \theta^{\text{old}})$

- For example, topic models.

• Solution: Restrict to distributions \mathcal{Q} that are easy to work with.

• The ELBO now looser: Intractable! We compute this

$$q^* = \operatorname{argmin}_{\mathcal{Q}} \text{KL}[q(\mathbf{z}), p(\mathbf{z} | \theta^{\text{old}})] \quad \text{for } \mathcal{Q} \in \mathcal{Q}$$

• Find an easy-to-work variational distribution q^* to approximate the inference distribution $p(\mathbf{z} | \mathbf{x}, \theta^{\text{old}})$.

- This group of methods are called Variational Methods.

$$\begin{aligned} Z &= \arg \max_{\mathbf{Z}} P(\mathbf{Z} | \mathbf{X}) \times \arg \max_{\mathbf{Z}} P(\mathbf{Z}, \mathbf{X}) \\ &= \arg \max_{\mathbf{Z}} \left[P(\mathbf{z}_1) \prod_{t=2}^T \left(\underbrace{\sum_{\mathbf{z}_t} P(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}_t)}_{\text{Transition}} \right) \cdot \underbrace{\sum_{\mathbf{z}_T} P(\mathbf{z}_T | \mathbf{x}_T)}_{\text{Emission}} \right] \\ &= \arg \min_{\mathbf{Z}} \left[\phi_1(\mathbf{z}_1) + \sum_{t=2}^T \left(\underbrace{\phi_t(\mathbf{z}_t) + \psi_t(\mathbf{z}_t, \mathbf{z}_{t-1})}_{\text{Cost}} \right) \right] \end{aligned}$$

minimizing negative log-probability.

是关于 (由于)

Lattice Diagram

$Z = \arg \min_{\mathbf{Z}} \left[\phi_1(\mathbf{z}_1) + \sum_{t=2}^T \left(\phi_t(\mathbf{z}_t) + \psi_t(\mathbf{z}_t, \mathbf{z}_{t-1}) \right) \right]$

• Node weight: $\phi_t(\mathbf{z}_t)$, edge weight $\psi_t(\mathbf{z}_t, \mathbf{z}_{t-1})$

• Row for each possible state, column for each time step.

• State transition diagram determines allowable paths (K^T in total).

• MAP estimate is the shortest path through weighted edges in the graph.

State I State 2 State K

$A_{11} \quad A_{12} \quad A_{13} \quad \dots$
 $A_{21} \quad A_{22} \quad A_{23} \quad \dots$
 $\vdots \quad \vdots \quad \vdots \quad \vdots$

• Viterbi Algorithm

$Z = \arg \min_{\mathbf{Z}} \left[\phi_1(\mathbf{z}_1) + \sum_{t=2}^T \left(\phi_t(\mathbf{z}_t) + \psi_t(\mathbf{z}_t, \mathbf{z}_{t-1}) \right) \right]$

• Define cost of shortest path ending with a specified state \mathbf{z}_T :

$$f_t(\mathbf{z}_t) = \phi_t(\mathbf{z}_t) + \min_{\mathbf{z}_{t-1}} \psi_t(\mathbf{z}_t, \mathbf{z}_{t-1}) + f_{t-1}(\mathbf{z}_{t-1})$$

• Interpretation: every path up to time t has the following form:

1. A path to some state at time $t-1$.
2. A cost for the transition from time $t-1$ to time t .

Initiation: $f_1(\mathbf{z}_1) = \phi_1(\mathbf{z}_1)$

At time step $t-1$, keep the shortest path and corresponding cost $f_{t-1}(\mathbf{z}_{t-1})$ ending with each state \mathbf{z}_{t-1} :

From time t to $t+1$:

- Evaluate the weights $\psi_t(\mathbf{z}_t, \mathbf{z}_{t-1})$ of every edge from t to $t+1$.

- K paths arriving at each node k , $O(K^2)$ computation cost in total.

Update the best path ending with each state \mathbf{z}_t .

$$f_t(\mathbf{z}_t) = \phi_t(\mathbf{z}_t) + \min_{\mathbf{z}_{t-1}} [\psi_t(\mathbf{z}_t, \mathbf{z}_{t-1}) + f_{t-1}(\mathbf{z}_{t-1})]$$

Algorithm stop:

- When reaching the end T , we have K paths through this lattice.

- Select the best one as the final choice.

Overall computation cost: $O(TK^2)$

Dynamical programming

7 伯努利混合模型的EM 算法 (15分)

考虑 $x = (x_1, \dots, x_D) \in \{0, 1\}^D$, 设 $P(x_d = 1) = p_d, 1 \leq d \leq D$, 令 $p = (p_1, \dots, p_D)$, 则 x 服从如下伯努利分布 $P(x | p) = \prod_{d=1}^D p_d^{x_d} (1 - p_d)^{1-x_d}$ 。

现考虑 K 个伯努利分布的混合, 混合分布参数为 $\mathbf{p} = \{p^{(1)}, \dots, p^{(K)}\}$, 对应权重为 $\boldsymbol{\pi} = \{\pi_1, \dots, \pi_K\}$ 。在混合分布中采样 $X = \{x^{(i)}\}_{i=1, \dots, n}$, 则有

$$P(x^{(i)} | \mathbf{p}, \boldsymbol{\pi}) = \sum_k \pi_k P(x^{(i)} | p^{(k)}) .$$

考虑 EM 算法中的 Expectation 步骤, 对于 $x^{(i)}$, 引入隐变量 $z^{(i)} \in \{0, 1\}^K$, 如果 $x^{(i)}$ 取自第 k 个伯努利分布, 则有 $z_k^{(i)} = 1$, 否则为 0。

(a) 请说明下面公式:

$$\begin{aligned} P(z^{(i)} | \boldsymbol{\pi}) &= \prod_{k=1}^K \pi_k^{z_k^{(i)}}, \\ P(x^{(i)} | z^{(i)}, \mathbf{p}, \boldsymbol{\pi}) &= \prod_{k=1}^K [P(x^{(i)} | p^{(k)})]^{z_k^{(i)}} . \end{aligned}$$

(b) 对应 $X = \{x^{(i)}\}_{i=1, \dots, n}$, 令 $Z = \{z^{(i)}\}_{i=1, \dots, n}$, 请基于(a)中结论, 试证明数据与隐变量的似然 $P(Z, X | \mathbf{p}, \boldsymbol{\pi})$ 为

$$P(Z, X | \mathbf{p}, \boldsymbol{\pi}) = \prod_{i=1}^n \left[\prod_{k=1}^K [P(x^{(i)} | p^{(k)})]^{z_k^{(i)}} \right] \left[\prod_{k=1}^K \pi_k^{z_k^{(i)}} \right] .$$

(c) 令 $\eta(z_k^{(i)}) = \mathbb{E}[z_k^{(i)} | x^{(i)}, \mathbf{p}, \boldsymbol{\pi}]$, 试证明

$$\eta(z_k^{(i)}) = \frac{\pi_k \prod_{d=1}^D (p_d^{(k)})^{x_d^{(i)}} (1 - p_d^{(k)})^{1-x_d^{(i)}}}{\sum_{j=1}^K \pi_j \prod_{d=1}^D (p_d^{(j)})^{x_d^{(i)}} (1 - p_d^{(j)})^{1-x_d^{(i)}}} .$$

(d) 令 $\tilde{\mathbf{p}}, \tilde{\boldsymbol{\pi}}$ 作为我们想要最大化的新参数, \mathbf{p} 和 $\boldsymbol{\pi}$ 来自上一次迭代, 基于(b)(c)中结论, 试证明 Maximization 步骤为最大化下面的期望:

$$\begin{aligned} &\mathbb{E}[\log P(Z, X | \tilde{\mathbf{p}}, \tilde{\boldsymbol{\pi}}) | X, \mathbf{p}, \boldsymbol{\pi}] \\ &= \sum_{i=1}^n \sum_{k=1}^K \eta(z_k^{(i)}) \left[\log \tilde{\pi}_k + \sum_{d=1}^D (x_d^{(i)} \log \tilde{p}_d^{(k)} + (1 - x_d^{(i)}) \log (1 - \tilde{p}_d^{(k)})) \right] \end{aligned}$$

(a) $x^{(i)}$ 必然恰好取 K 个伯努利分布的其中一个, 由于这 K 个分布对应权重分别为 $\{\pi_1, \dots, \pi_K\}$, 取在第 k 个的概率为 π_k 。此时只有

$$z_{k_0}^{(i)} = 1, \text{ 其余 } K-1 \text{ 维上均为 } 0, \text{ 作为指数时幂的取值为 } 1, \therefore P(z^{(i)} | \boldsymbol{\pi}) = \pi_{k_0} = \prod_{k=1}^K \pi_k^{z_k^{(i)}} .$$

同理, $P(x^{(i)} | z^{(i)}, \mathbf{p}, \boldsymbol{\pi}) = P(x^{(i)} | p^{(k_0)}) = \prod_{k=1}^K [P(x^{(i)} | p^{(k)})]$, 即 $z^{(i)}, \mathbf{p}, \boldsymbol{\pi}$ 的条件下满足 $x = x^{(i)}$ 的 x 必然来自某个具体的第 k_0 个伯努利分布。

$$(b) \text{ 通过似然的表示及简单代入可得, } P(Z, X | \mathbf{p}, \boldsymbol{\pi}) = p(z^{(i)} | \boldsymbol{\pi}) \prod_{i=1}^n P(x^{(i)} | z^{(i)}, \mathbf{p}, \boldsymbol{\pi}) = \prod_{k=1}^K [\prod_{k=1}^K [P(x^{(i)} | p^{(k)})]^{z_k^{(i)}}] [\prod_{k=1}^K \pi_k^{z_k^{(i)}}] .$$

$$\textcircled{(c)} \text{ 注意到 } z_k^{(i)} = k \text{ 当且仅当 } x^{(i)} \text{ 取自第 } k \text{ 个伯努利分布, 故 } \eta(z_k^{(i)}) = E[z_k^{(i)} | x^{(i)}, \mathbf{p}, \boldsymbol{\pi}] = p(k | x^{(i)}, \mathbf{p}, \boldsymbol{\pi}) = \frac{\pi_k p(x^{(i)} | p^{(k)})}{\sum_{j=1}^K \pi_j p(x^{(i)} | p^{(j)})} , \text{ 由伯努利分}$$

$$\text{布的形式, 有 } p(x^{(i)} | p^{(k)}) = \prod_{d=1}^D ((p_d^{(k)})^{x_d^{(i)}} (1 - p_d^{(k)})^{1-x_d^{(i)}}), \text{ 代入期望表达式中即得 } E[z_k^{(i)} | x^{(i)}, \mathbf{p}, \boldsymbol{\pi}] = \frac{\pi_k \prod_{d=1}^D ((p_d^{(k)})^{x_d^{(i)}} (1 - p_d^{(k)})^{1-x_d^{(i)}})}{\sum_{j=1}^K \pi_j \prod_{d=1}^D ((p_d^{(j)})^{x_d^{(i)}} (1 - p_d^{(j)})^{1-x_d^{(i)}})} .$$

$$\text{(d) 对(b)中的 } P(Z, X | \mathbf{p}, \boldsymbol{\pi}) \text{ 取对数, 有 } \log P(Z, X | \mathbf{p}, \boldsymbol{\pi}) = \sum_{i=1}^n \log \prod_{k=1}^K [P(x^{(i)} | p^{(k)})]^{z_k^{(i)}} + \sum_{k=1}^K z_k^{(i)} \log \pi_k = \sum_{k=1}^K z_k^{(i)} (\log \pi_k + \sum_{i=1}^n \log P(x^{(i)} | p^{(k)})), \text{ 其中 } \log P(x^{(i)} | p^{(k)}) = \sum_{d=1}^D (x_d \log p_d + (1 - x_d) \log (1 - p_d)) , \text{ 恰与结果中后一方括号内求和的通项一致, 而整体求和中的 } z_k^{(i)} \text{ 项使得和恰好等于其中使得 } x^{(i)} \text{ 取第 } k \text{ 个伯努利分布的第 } k \text{ 项, 故 } \log P(Z, X | \mathbf{p}, \boldsymbol{\pi}) \text{ 与结果中方括号内的形式一致 (只是 } p \text{ 和 } \boldsymbol{\pi} \text{ 为个别取值而非均值), 而 } E[\log P(Z, X | \tilde{\mathbf{p}}, \tilde{\boldsymbol{\pi}}) | X, \mathbf{p}, \boldsymbol{\pi}] = \sum_{j=1}^K \eta(z_k^{(i)}) \log P(Z, X | \mathbf{p}, \boldsymbol{\pi}) . \text{ 故}$$

$$E[\log P(Z, X | \tilde{\mathbf{p}}, \tilde{\boldsymbol{\pi}}) | X, \mathbf{p}, \boldsymbol{\pi}] = \sum_{i=1}^n \sum_{k=1}^K \eta(z_k^{(i)}) [\log \tilde{\pi}_k + \sum_{d=1}^D (x_d^{(i)} \log \tilde{p}_d^{(k)} + (1 - x_d^{(i)}) \log (1 - \tilde{p}_d^{(k)}))] , \text{ 即求得该形式。}$$