

Machine Learning 2

Information Theory

- Relative Entropy**
- revisited**

Patrick Forré

Relative Entropy / Kullback-Leibler Divergence

- Let $q(x)$ and $p(x)$ be two distribution for the same observation space.
- Then the Relative Entropy or Kullback-Leibler Divergence is defined:

$$\text{KL} (q(X) \| p(X)) := \mathbb{E}_{q(X)} \left[\log \frac{q(X)}{p(X)} \right]$$

- Fundamental Inequality:

$$\text{KL} (q(X) \| p(X)) \geq 0,$$

with equality if and only if $q(x) = p(x)$ for almost all values x .

Chain Rule and Data Processing Inequality for KL

- Let $q(x, z)$ and $p(x, z)$ two joint distribution for the same product space.
- Then we have the Chain Rule:

$$\text{KL}(q(X, Z) \| p(X, Z)) = \underset{\geq 0}{\text{KL}}(q(x) \| p(x)) + \underset{\geq 0}{\mathbb{E}_{q(x)}}[\text{KL}(q(z|x) \| p(z|x))]$$

- From this we get the Data Processing Inequality:

$$\text{KL}(q(X, Z) \| p(X, Z)) \geq \text{KL}(q(x) \| p(x))$$

with equality if and only if $q(z|x) = p(z|x)$ for almost all values x and all z .

Proof: $\log \frac{q(z|x) q(x)}{p(z|x) p(x)} = \log \frac{q(z|x)}{p(z|x)} + \log \frac{q(x)}{p(x)}$

Principle of Minimal Relative Entropy

- Let $q(x)$ be an underlying ‘true’ distribution of interest and $\{p_\theta(x) \mid \theta \in \Theta\}$ a statistical model for $q(x)$.
- The Principle of Minimal Relative Entropy says one should try to minimize:

$$\theta^* \in \arg \min_{\theta} \text{KL} (q(X) \| p_{\theta}(X))$$

- We have seen that if one replaces the (unknown) $q(x)$ with an empirical distribution $\hat{q}(x)$ one arrives at the learning framework of Maximum Likelihood Estimation.

Machine Learning 2

**Variational Inference
- Ideas and Principles**

Patrick Forré

Remark: Approximate Inference

- If $q(x)$ is the distribution / marginal of interest and d a divergence measure quantifying the difference between $q(x)$ and a model distribution $p_\theta(x)$,

- then minimizing such a divergence (e.g. $d = \text{KL}$):

$$\arg \min_{\theta} d(q(X) \| p_{\theta}(X))$$

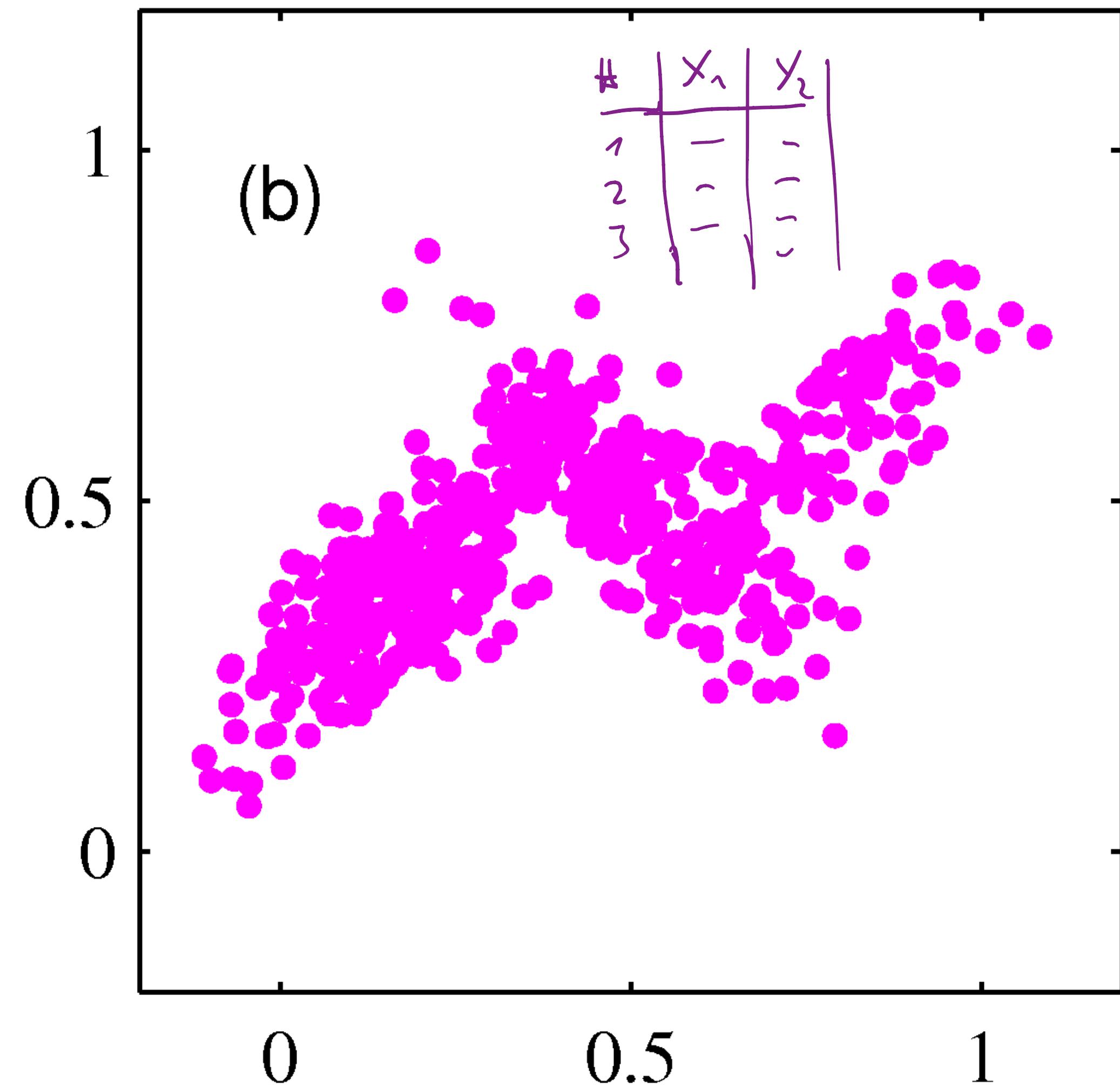
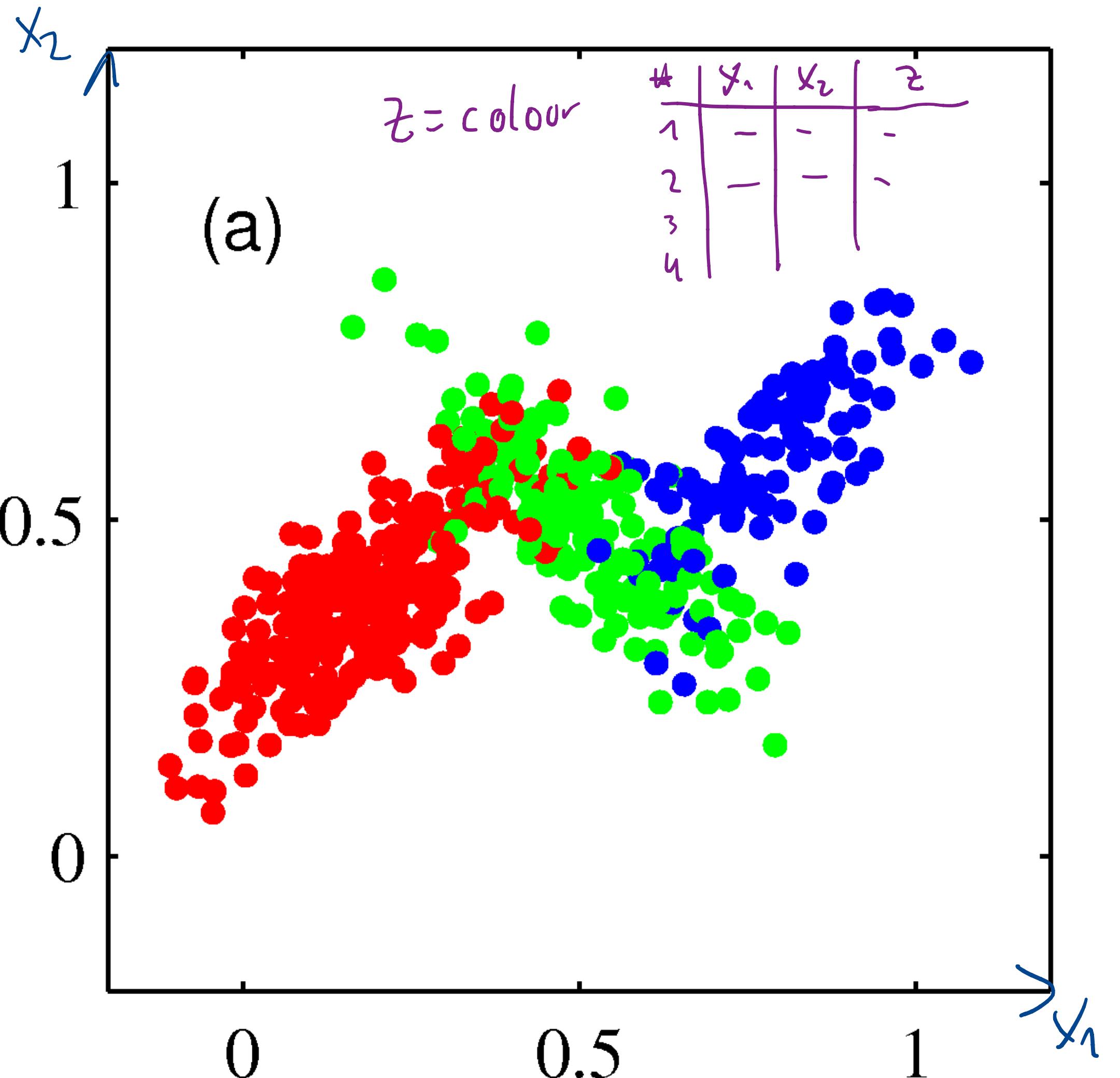
between $q(x)$ and an approximating model $p_\theta(x)$ can be seen as a general approach for approximate inference.

- Small values for $d(q(X) \| p_{\theta}(X))$ means good approximation.

Remark: Latent Variable Model vs Bayesian Setting

- Consider joint distribution $p_{\theta}(x, z) = p_{\theta_2}(x | z) \cdot p_{\theta_1}(z)$ and data \hat{X} .
- Latent Variable Model
 - z : latent variables
 - θ : parameters
 - $p_{\theta_1}(z)$: marginal latent distribution
 - $p_{\theta}(\hat{X})$: likelihood
 - $p_{\theta_2}(x | z)$: latent conditional distribution
 - $p_{\theta}(z | x)$: conditional distribution
 - Fitting θ to \hat{X} : Maximum Likelihood Estimation
 - Names and interpretations blur, but the same math applies!
- Bayesian Setting
 - z : parameters
 - θ : hyper-parameters
 - $p_{\theta_1}(z)$: prior distributions
 - $p_{\theta}(\hat{X})$: evidence
 - $p_{\theta_2}(x | z)$: likelihood
 - $p_{\theta}(z | x)$: posterior distribution
 - Fitting θ to \hat{X} : Empirical Bayes (max. evidence)

Example - Latent Variable Model - Clustering



Two approaches: MLE + VI

- Let $q(x)$ be the ‘true’ distribution of interest and $\{p_\theta(x, z) = p_{\theta_2}(x | z) \cdot p_{\theta_1}(z) \mid \theta \in \Theta\}$ a (generative) **latent variable** model for $q(x)$.
- To compare $q(x)$ to $p_\theta(x, z)$ we have **two** approaches:
 - 1.) • **marginalize** $p_\theta(x, z)$ to $p_\theta(x)$ and minimize: $\min_{\theta} \text{KL} (q(X) \parallel p_\theta(X)) \rightarrow \text{MLE}$
 - 2.) • introduce a **variational** family of **inference** distributions: $\{q_w(z | x) \mid w \in \mathcal{W}\} \rightarrow \text{VI}$ (meant to approximate $p_\theta(z | x)$) and minimize:
$$\min_{w, \theta} \text{KL} (q_w(X, Z) \parallel p_\theta(X, Z))$$
- Replacing expectation $\mathbb{E}_{q(X)}$ with an **empirical mean** $\mathbb{E}_{\hat{q}(X)}$ leads in 1st case to MLE and in the 2nd case to the learning framework of **Variational Inference**.

Maximum Likelihood Estimation vs Variational Inference

- By the **chain rule:**

$$\mathbb{E}_{q(X)} \left[\text{KL} \left(q_w(Z|X) \| p_\theta(Z|X) \right) \right] + \text{KL} \left(q(X) \| p_\theta(X) \right) = \text{KL} \left(q_w(X, Z) \| p_\theta(X, Z) \right)$$

MLE *VI*
VI VI
O O

we see that Variational Inference degenerates to MLE if the inference model is flexible enough to minimize the left hand term.

- We have the inequality for w, θ :

$$0 \leq \text{KL} \left(q(X) \| p_\theta(X) \right) \leq \text{KL} \left(q_w(X, Z) \| p_\theta(X, Z) \right)$$

with (rhs) **equality** iff $q_w(z|x) = p_\theta(z|x)$ for all values x, z .

Empirical versions of the KL divergencies

- $$-\text{KL}(q(X) \| p_\theta(X)) - H(q(X)) = -\mathbb{E}_{q(x)} \left[\log \frac{q(x)}{p_\theta(x)} \right] - \mathbb{E}_{q(x)} [-\log q(x)]$$

$$= \mathbb{E}_{q(x)} [\log p_\theta(x)]$$

$\stackrel{\text{LLN}}{\approx} \mathbb{E}_{q(x)} [\log p_\theta(x)] = \frac{1}{n} \log p_\theta(\hat{x})$

evidence / log-likelihood

VI
- $$-\text{KL}(q_w(X, Z) \| p_\theta(X, Z)) - H(q(X)) \approx \frac{1}{n} \mathbb{E}_{q_w(z|x)} \left[\log \frac{p_\theta(z|x)}{q_w(z|x)} \right]$$

evidence - lower bound
(ELBO)

$$= -\mathbb{E}_{q_w(z|x)q(x)} \left[\log \frac{q_w(z|x)q(x)}{p_\theta(z|x)} \right] - \mathbb{E}_{q(x)} [-\log q(x)]$$

$$= +\mathbb{E}_{q_w(z|x)q(x)} \left[\log \frac{p_\theta(z|x)}{q_w(z|x)} \right] - \mathbb{E}_{q_w(z|x)q(x)} [\log q(x)] + \mathbb{E}_{q(x)} [\log q(x)]$$

$$= \mathbb{E}_{q(x)} \mathbb{E}_{q_w(z|x)} \left[\log \frac{p_\theta(z|x)}{q_w(z|x)} \right] \stackrel{\text{LLN}}{\approx} \mathbb{E}_{\hat{q}(x)} \mathbb{E}_{q_w(z|x)} \left[\log \frac{p_\theta(z|x)}{q_w(z|x)} \right] = \frac{1}{N} \mathbb{E}_{q_w(z|x)} \left[\log \frac{p_\theta(z|x)}{q_w(z|x)} \right]$$

Empirical versions: Log-Likelihood and Evidence Lower Bound

- Log-Likelihood: $L_{\hat{X}}(\theta) = \log p_\theta(\hat{X}) =$
- ELBO: $L_{\hat{X}}(w, \theta) =$

Fundamental Lemma of Variational Inference

- Let $p_\theta(z, x)$ be a latent variable model and \hat{X} data and $q_w(z | x)$ be an inference model (meant to approximate $p_\theta(z | \hat{X})$).
- Then the evidence lower bound is a lower bound to the evidence (aka log-likelihood):

$$\log p_\theta(\hat{X}) = L_{\hat{X}}(\theta) \geq L_{\hat{X}}(w, \theta) = \mathbb{E}_{q_w(Z|\hat{X})} \left[\log \frac{p_\theta(Z, \hat{X})}{q_w(Z|\hat{X})} \right]$$

ELBO

- Equality holds iff $q_w(z | \hat{X}) = p_\theta(z | \hat{X})$ for all z .
- Proof: Arguments from before or directly with Jensen's inequality.

Proof

$$\log P(\hat{x}) = \log \int \frac{p(z, \hat{x})}{q(z|\hat{x})} \cdot q(z|\hat{x}) dz$$

$$= \log \mathbb{E}_{q(z|\hat{x})} \left[\frac{p(z, \hat{x})}{q(z|\hat{x})} \right]$$

$$\geq \mathbb{E}_{q(z|\hat{x})} \left[\log \frac{p(z, \hat{x})}{q(z|\hat{x})} \right] = ELBO$$

□

The Framework of Variational Inference (VI)

- Let $q(x)$ be an underlying ‘true’ distribution of interest and \hat{X} data sampled from it.
Let $\{p_\theta(x, z) \mid \theta \in \Theta\}$ be a **latent variable model** for $q(x)$ (or a Bayesian setting).
- Instead of **maximizing the evidence/log-likelihood** $\log p_\theta(\hat{X})$ w.r.t. θ ,
- in **Variational Inference (VI)** one introduces a variational family of **inference distributions** $\{q_w(z \mid x) \mid w \in \mathcal{W}\}$ and **Maximizes the Evidence Lower BOund (ELBO)**:

$$L_{\hat{X}}(w, \theta) := \mathbb{E}_{q_w(Z \mid \hat{X})} \left[\log \frac{p_\theta(Z, \hat{X})}{q_w(Z \mid \hat{X})} \right]$$

w.r.t. **both parameters** w, θ .

- VI was justified either by the Minimal Relative Entropy Principle applied to the joint distributions plus the empirical approximation of $\mathbb{E}_{q(X)}$ via the law of large numbers, or via the Maximum Likelihood Estimation and its (log) approximation by the ELBO from below.

Machine Learning 2

**Variational Inference
- Expectation-Maximization Algorithm**

Patrick Forré

Expectation-Maximization Algorithm

- Let $q(x)$ be an underlying ‘true’ distribution of interest and \hat{X} data sampled from it. Let $\{p_\theta(x, z) \mid \theta \in \Theta\}$ be a **latent variable** model for $q(x)$ and $\{q_w(z|x) \mid w \in \mathcal{W}\}$ a variational **family of inference distributions**.
- To **maximize the ELBO** (VI): $L_{\hat{X}}(w, \theta) := \mathbb{E}_{q_w(Z|\hat{X})} \left[\log \frac{p_\theta(Z, \hat{X})}{q_w(Z|\hat{X})} \right]$ we **alternate**:
 - **E-step:** $w^{(t+1)} := \underset{w}{\operatorname{argmax}} \ L_{\hat{X}}(w, \theta^{(t)}) + \lambda \cdot (\text{constraints})$
 - **M-step:** $\theta^{(t+1)} := \underset{\theta}{\operatorname{argmax}} \ L_{\hat{X}}(w^{(t+1)}, \theta) + \tilde{\lambda} \cdot (\text{constraints})$
- In practice we also need to include **Lagrange multiplier** in each optimization step to ensure constraints (e.g. normalization of probabilities).
- Note that the ELBO can only **improve** in every step, so the EM-algorithm converges.

Simplifying the steps

- **E-step:**

$$\begin{aligned} w^{(t+1)} &:= \underset{w}{\operatorname{argmax}} L_{\hat{x}}(w, \theta^{(t)}) = \underset{w}{\operatorname{argmax}} \mathbb{E}_{q_w(z|\hat{x})} \left[\log \frac{p_{\theta^{(t)}}(z, \hat{x})}{q_w(z|\hat{x})} \right] \\ &= \underset{w}{\operatorname{argmin}} \text{KL}\{q_w(z|\hat{x}) \parallel p_{\theta^{(t)}}(z|\hat{x}) \} \end{aligned}$$

- **M-step:**

$$\begin{aligned} \theta^{(t+1)} &:= \underset{\theta}{\operatorname{argmax}} L_{\hat{x}}(w^{(t+1)}, \theta) = \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{Q_{w^{(t+1)}}(z|x)} \left[\log \frac{p_{\theta}(z, \hat{x})}{q_{w^{(t+1)}}(z|\hat{x})} \right] \\ &= \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{q_{w^{(t+1)}}(z|\hat{x})} \left[\log p_{\theta}(z, \hat{x}) \right] \end{aligned}$$

The EM-Algorithm with best possible E-steps

- **E-step:**

$$\bullet w^{(t+1)} := \underset{w}{\operatorname{argmin}} \text{KL}(q_w(z|\hat{x}) \| p_{\theta^{(t)}}(z|\hat{x}))$$

- **M-step:**

$$\theta^{(t+1)} := \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{q_w^{(t+1)}(z|\hat{x})} [\log p_{\theta}(z|\hat{x})]$$

- The best possible **E-step** occurs when we can arrange $w^{(t+1)}$ such that:

$$\bullet q_{w^{(t+1)}}(Z|\hat{X}) = p_{\theta^{(t)}}(Z|\hat{X})$$

implying:

$$\bullet \mathbb{E}_{q_{w^{(t+1)}}(Z|\hat{X})} [\log p_{\theta}(Z|\hat{X})] = \mathbb{E}_{p_{\theta^{(t)}}(Z|\hat{X})} [\log p_{\theta}(Z|\hat{X})] =: Q(\theta, \theta^{(t)})$$

- So the objective in the **M-step** after a best possible E-step would be:

$$\bullet \theta^{(t+1)} := \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^{(t)})$$

The EM-Algorithm with best possible E-steps

- The EM-Algorithm in this case reduces to alternating:

- **E-step:** Evaluate: $Q(\theta, \theta^{(t)}) := \mathbb{E}_{P_{\theta^{(t)}}(z|x)} [\log P_\theta(z, \hat{x})]$
- **M-step:** Compute: $\theta^{(t+1)} := \arg \max_{\theta} Q(\theta, \theta^{(t)}) + \lambda(\text{constraint})$
- In the **M-step** don't forget **Lagrange multipliers** if there are **constraints**.
- The **M-step** can sometimes be solved by putting gradients to 0 or using gradient methods for optimization.

Machine Learning 2

Variational Inference

- Expectation-Maximization Algorithm**
- Example: Mixtures of Exponential Families**

Patrick Forré

$q(x)$ true distr. $x^{(1)}, \dots, x^{(N)}$, iid sample

$$P(z|\pi) = \prod_{k=1}^K \pi_k^{z_k}$$

$$z_k = \mathbb{1}_k(z)$$

$$z \in \{1, \dots, K\}$$

$$z_k^{(n)} := \mathbb{1}_k(z^{(n)})$$

$$z^{(n)} \hookrightarrow x^{(n)}$$

$$\hat{x} = [x^{(1)}, \dots, x^{(N)}]^T$$

$$P(x|z, y) = \prod_{k=1}^K \left(\exp(\gamma_k^T T_k(x) - A_k(y_k) + B_k(x)) \right)^{z_k}$$

$$z = (z_k^{(n)})_{\substack{n=1, \dots, N \\ k=1, \dots, K}}$$

$$P(z, \hat{x} | \pi, y) \stackrel{\text{iid}}{=} \prod_{n=1}^N \prod_{k=1}^K \exp \left[z_k^{(n)} \cdot (\log \pi_k + \gamma_k^T T_k(x^{(n)}) + B_k(x^{(n)}) - A_k(y_n)) \right]$$

$$\log P(z, \hat{x} | \pi, y) = \sum_{n,k} z_k^{(n)} \cdot [\gamma_k^T T_k(x^{(n)}) + B_k(x^{(n)}) - A_k(y_n) + \log \pi_k]$$

$$\log p(z_i | \hat{x}_i | \pi_i, y) = \sum_{k \in K} z_k^{(n)} \cdot \left[\gamma_k^T T_k(x^{(n)}) + B_k(x^{(n)}) - A_k(y_k) + \log \pi_k \right]$$

E-step $\mathbb{E}_{p(z | \hat{x}, \tilde{\pi}, \tilde{\gamma})} [\log p(z_i | \pi_i, y)] =: Q((\pi_i, y), (\tilde{\pi}, \tilde{\gamma}))$

$$j_k(x^{(n)}) := \mathbb{E}_{p(z | \hat{x}, \tilde{\pi}, \tilde{\gamma})} [z_k^{(n)}] = p(z^{(n)} = k | x^{(n)}, \tilde{\gamma}, \tilde{\pi})$$

$$p(z^{(n)}, x^{(n)} | \tilde{\pi}, \tilde{\gamma}) = \prod_{q=1}^K \tilde{\pi}_q^{z_q^{(n)}} \cdot \exp(\tilde{\gamma}_q^T T_q(x^{(n)}) + B_q(x^{(n)}) - A_q(y_q))^{z_q^{(n)}}$$

$$p(x^{(n)} | \tilde{\pi}, \tilde{\gamma}) = \sum_{q=1}^K \tilde{\pi}_q \cdot \exp(\tilde{\gamma}_q^T T_q(x^{(n)}) + B_q(x^{(n)}) - A_q(y_q))$$

$$\rightarrow j_q(x^{(n)}) = \frac{\tilde{\pi}_q \cdot \exp(\tilde{\gamma}_q^T T_q(x^{(n)}) + B_q(x^{(n)}) - A_q(y_q))}{\sum_{j=1}^K \tilde{\pi}_j \cdot \exp(\tilde{\gamma}_j^T T_j(x^{(n)}) + B_j(x^{(n)}) - A_j(y_j))}$$

E-step :

$$Q((\pi_1, \gamma), (\tilde{\pi}_1, \tilde{\gamma})) := \sum_{k, \ell} \gamma_k(x^{(n)}) \cdot [\gamma_k^\top T_\ell(x^{(n)}) + B_\ell(x^{(n)}) - A_k(\gamma_k) + \log \tilde{\pi}_k]$$

where

$$\gamma_k(x^{(n)}) := \frac{\tilde{\pi}_k \cdot \exp(\tilde{\gamma}_k^\top T_k(x^{(n)}) + B_k(x^{(n)}) - A_k(\tilde{\gamma}_k))}{\sum_{j=1}^K \tilde{\pi}_j \cdot \exp(\tilde{\gamma}_j^\top T_j(x^{(n)}) + B_j(x^{(n)}) - A_j(\tilde{\gamma}_j))}$$

$$\text{M-step : } (\hat{\pi}_l^*, \hat{y}^*) = \arg \max_{\pi_l, y} Q((\pi_l, y), (\hat{\pi}, \hat{y}))$$

$$(\hat{\pi}^*, \hat{y}^*) = \arg \max_{\pi_l, y} \sum_{n, k} y_k x^{(n)} \cdot [\gamma_k^\top T_k(x^{(n)}) + B_k(x^{(n)}) - A_k(y_k) + \log \pi_k]$$

$$0 \stackrel{?}{=} \frac{\partial Q}{\partial y_n} = \sum_{k=1}^N y_k x^{(n)} \cdot [T_k(x^{(n)}) - \underbrace{A_k^t(y_k)}_{\mathbb{E}[T_k(x)]|y_k}]$$

$$\rightarrow A_k^t(y_k) = \frac{1}{N} \sum_{n=1}^N y_k x^{(n)} \cdot T_k(x^{(n)})$$

$$\rightarrow \gamma_k = (A_k^t)^{-1} \left(\frac{1}{N} \sum_{n=1}^N y_k(x^{(n)}) \cdot T_k(x^{(n)}) \right)$$

(Or $\gamma_k := \gamma_k + \alpha \cdot \left(\sum_{n=1}^N y_k(x^{(n)}) [T_k(x^{(n)}) - A_k^t(y_k)] \right) \underbrace{\mathbb{E}[T_k(x)|y_k]}_{\text{or use samples to approximate}}$)

M-step:

$$(\pi^*, \gamma^*) = \underset{\pi, \gamma}{\operatorname{argmax}}$$

$$\underbrace{Q((y, \pi), (\tilde{y}, \tilde{\pi}))}_{\sum_{k, n} \gamma_k(x^{(n)}) \cdot [\gamma_k^\top \tau_k(x^{(n)}) + B_k(x^{(n)}) - A_k(y_k) + \log \pi_k]}$$

constraints: $\sum_{k=1}^K \pi_k = 1 \quad \sim \quad \mathcal{L}(y, \pi) := Q(y, \pi) + \lambda \cdot \left(1 - \sum_{k=1}^K \pi_k\right)$

$$0 \stackrel{?}{=} \frac{\partial \mathcal{L}}{\partial \pi_k} = \left(\sum_{n=1}^N \frac{1}{\pi_k} \gamma_k(x^{(n)}) \right) - \lambda$$

$$\pi_k = \frac{1}{N} \sum_{n=1}^N \gamma_k(x^{(n)})$$

$$\Rightarrow \lambda \cdot \pi_k = \sum_{n=1}^N \gamma_k(x^{(n)})$$

$$\Rightarrow \lambda = \lambda \cdot 1 = \sum_{k=1}^K \lambda \cdot \pi_k = \sum_{k, n} \gamma_k(x^{(n)}) = \sum_{n=1}^N \underbrace{\sum_k \gamma_k(x^{(n)})}_1 = N$$

E-step : $\forall \alpha=1\dots K, \forall \gamma=1\dots N:$

$$\gamma_{\alpha}^{(t+1)}(x^{(\gamma)}) := \frac{\pi_{\alpha}^{(t)} \cdot \exp\left(\eta_{\alpha}^{(t)} + T_{\alpha}(x^{(\gamma)}) + B_{\alpha}(x^{(\gamma)}) - A_{\alpha}(\gamma_{\alpha}^{(t)})\right)}{\sum_{j=1}^K \pi_j^{(t)} \cdot \exp\left(\eta_j^{(t)} + T_j(x^{(\gamma)}) + B_j(x^{(\gamma)}) - A_j(\gamma_j^{(t)})\right)}$$

M-step $\forall \alpha=1\dots K:$

$$\pi_{\alpha}^{(t+1)} := \frac{1}{N} \sum_{n=1}^N \gamma_{\alpha}^{(t+1)}(x^{(\gamma)})$$

$$\eta_{\alpha}^{(t+1)} := (A_{\alpha}^t)^{-1} \left(\frac{1}{N} \sum_{n=1}^N \gamma_{\alpha}^{(t+1)}(x^{(\gamma)}) \cdot T_{\alpha}(x^{(\gamma)}) \right)$$