



• Entropy of Exponential Family and Graphical Model

$$H_p(x) = -E_{x \sim p} [\ln p(x)]$$



low entropy: few low-prob instances

$$= -E_{x \sim p} \left[\ln \frac{1}{z(t(\theta))} \cdot \exp\langle t(\theta), \eta(x) \rangle \right]$$

$$= \ln z(t(\theta)) - \langle t(\theta), E_{x \sim p}[\eta(x)] \rangle$$



g. entropy of Gaussian $X \sim N(\mu, \sigma^2)$

$$t(\theta) \cdot \eta(\theta) = -\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2} = -\frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} - \frac{\mu^2}{2\sigma^2}$$

$$t(\theta) = \begin{bmatrix} -\frac{1}{2\sigma^2} \\ \frac{\mu}{\sigma^2} \\ 1 \end{bmatrix}, \eta(x) = \begin{bmatrix} x^2 \\ x \\ 1 \end{bmatrix}$$

$$\begin{aligned} H_p(x) &= \ln \frac{1}{z(t(\theta))} - \left(\frac{1}{2} E[x^2] + \frac{\mu}{\sigma^2} E[x] - \frac{\mu^2}{2\sigma^2} \right) \quad E[X] = \mu + \sigma^2 \\ &= \frac{1}{2} \ln 2\pi\sigma^2 + \frac{(\mu+\sigma^2)^2}{2\sigma^2} - \frac{\mu^2}{2\sigma^2} + \frac{\mu^2}{2\sigma^2} \\ &= \frac{1}{2} \ln 2\pi\sigma^2 + \frac{1}{2} = \ln \sqrt{2\pi\sigma^2} \end{aligned}$$

if $p(x) = \frac{1}{Z} T_k \phi_k(D)$ is a markov network, then $H_p(x) = \ln Z + \sum_i E_p[\ln \phi_i(D)]$
 $\ln Z = \ln T_k \phi_k(D)$

if $P(X) = T_p(p(x_i | p_{\theta}(x_i)))$ is a distribution over a Bayesian network, then

$$H_p(x) = \sum_i H_p(x_i | p_{\theta}(x_i))$$

$$H_p(x) = E_p[-\ln p(x)] = E_p[-\sum_i \ln p(x_i | p_{\theta}(x_i))] = \sum_i E_p[-\ln p(x_i | p_{\theta}(x_i))] = \sum_i H_p(x_i | p_{\theta}(x_i))$$

• Relative Entropy

$$D(P||Q) = E_{x \sim P} \left[\ln \frac{p(x)}{q(x)} \right]$$

if Q is a distribution and P_θ is an exponential family defined by $\eta(\cdot)$ and $t(\cdot)$

$$D(Q || P_\theta) = E_{x \sim Q} \left[\ln \frac{p(x)}{p_\theta(x)} \right] = E_{x \sim Q} [\ln q(x)] + E_{x \sim Q} \left[\ln \frac{1}{z(t(\theta))} \cdot \exp\langle t(\theta), \eta(x) \rangle \right]$$

$$= -H_q(x) + \ln z(t(\theta)) - \langle t(\theta), E_{x \sim Q}[\eta(x)] \rangle$$

if P_θ and $P_{\theta'}$ are in the same exponential family

$$D(P_\theta || P_{\theta'}) = E_{x \sim P_\theta} \left[\ln \frac{p_\theta(x)}{p_{\theta'}(x)} \right]$$

$$= E_{x \sim P_\theta} \left[\ln \frac{1}{z(t(\theta))} \exp\langle t(\theta), \eta(x) \rangle \right] - E_{x \sim P_\theta} \left[\ln \frac{1}{z(t(\theta'))} \cdot \exp\langle t(\theta'), \eta(x) \rangle \right]$$

$$= -\ln \frac{z(t(\theta))}{z(t(\theta'))} + E_{x \sim P_\theta} [\langle t(\theta), \eta(x) \rangle - \langle t(\theta'), \eta(x) \rangle]$$

$$= -\ln \frac{z(t(\theta))}{z(t(\theta'))} + \langle E_{x \sim P_\theta} [\eta(x)], t(\theta) - t(\theta') \rangle$$

If P is a distribution over a Bayesian Network G ,

$$D(Q||P) = \mathbb{E}_{x \sim P} [\ln Q] - \mathbb{E}_{x \sim P} [\ln P]$$

$$= -H_Q(X) - \sum_{x \sim P(x)} E_{x \sim P(x)} [\ln P(x|f(x))]$$

$$\begin{aligned} \mathbb{E}_{x \sim P(x)} [\ln P(x|f(x))] &= \int_{\Omega} \int_{\Omega^X} Q(x_i, f(x_i)) \cdot \ln p(x_i|f(x_i)) dx_i d\Omega \\ &= \int_{\Omega} \int_{\Omega^X} Q(f(x_i)) \cdot Q(x_i|f(x_i)) \cdot \ln p(x_i|f(x_i)) dx_i d\Omega \\ &= \int_{\Omega} Q(f(x_i)) \sum_{x \sim P(x)} E_{x \sim P(x|f(x))} [\ln P(x_i|f(x_i))] \\ &= -H_Q(X) - \sum_{x \sim P(x)} \sum_{f(x)=y} Q(f(x)) E_{x \sim P(x|f(x))} [\ln P(x_i|f(x_i))] \end{aligned}$$

$$\begin{aligned} &\int_{\Omega} \int_{\Omega^X} P(x_i, x_k) \ln P(x_i|f(x_i)) dx_i d\Omega \\ &= \int_{\Omega} \int_{\Omega^X} P(x_i, x_k) \ln p(x_i|f(x_i)) dx_i d\Omega + \int_{\Omega} \int_{\Omega^X} P(x_i, x_k) \ln f(x_i) dx_i d\Omega \\ &= \int_{\Omega} P(x_i) \ln f(x_i) dx_i + \dots \\ &= \mathbb{E}_{x \sim P} [\ln f(x)] + \mathbb{E}_{x \sim P} [\ln P(x_i|f(x_i))] \end{aligned}$$

If P, Q are 2 distributions over G ,

$$D(Q||P) = \mathbb{E}_{x \sim P} [\ln Q] - \mathbb{E}_{x \sim P} [\ln P]$$

$$\begin{aligned} &= \sum_i \sum_{f(x_i)} Q(f(x_i)) \mathbb{E}_{x \sim P(x|f(x))} [\ln Q(x_i|f(x_i))] - \sum_i \sum_{f(x_i)} Q(f(x_i)) \mathbb{E}_{x \sim P(x|f(x))} [\ln P(x_i|f(x_i))] \\ &= \sum_i \sum_{f(x_i)} Q(f(x_i)) D(Q(x_i|f(x_i)) \| P(x_i|f(x_i))) \end{aligned}$$

Projection

View relative entropy as a notion of distance between two distributions

Let P be a distribution and let \mathcal{Q} be a convex set of distribution

I-projection (information projection) of P onto \mathcal{Q} : $\mathcal{Q}^I = \arg \min_{Q \in \mathcal{Q}} D(Q||P)$

M-projection (moment projection) of P onto \mathcal{Q} : $\mathcal{Q}^M = \arg \min_{Q \in \mathcal{Q}} D(P||Q)$

(if $P \in \mathcal{Q}$, then $\mathcal{Q}^I = \mathcal{Q}^M = P$)

because the relative entropy is not symmetric, $\mathcal{Q}^M \neq \mathcal{Q}^I$ in general

M-projection tends to give all assignments high density; I-projection tend to have some instances with low density

e.g. project P onto Gaussians

$$I: D(P||Q) = \int_X P(x) \ln \frac{P(x)}{Q(x)} dx$$

$$= -H_P(X) - \mathbb{E}_{x \sim P} [\ln Q(x)]$$



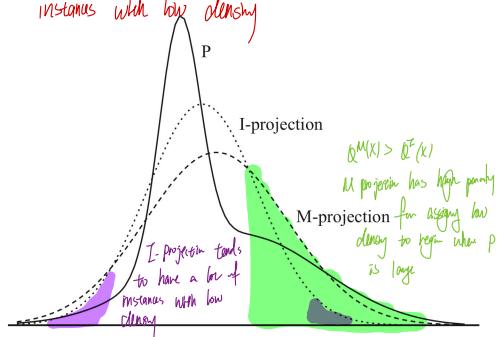
$$Q^M = \arg \max_{Q \in \mathcal{Q}} -H_Q(X) - \mathbb{E}_{x \sim P} [\ln Q(x)] = \arg \max_{Q \in \mathcal{Q}} \mathbb{E}_{x \sim P} [\ln Q(x)]$$

large penalty for assigning low density to regions where P is large.

$$I: D(Q||P) = \int_X Q(x) \ln \frac{Q(x)}{P(x)} dx$$

$$= -H_Q(X) - \mathbb{E}_{x \sim Q} [\ln P(x)]$$

$$Q^I = \arg \min_{Q \in \mathcal{Q}} -H_Q(X) + \mathbb{E}_{x \sim Q} [\ln P(x)] = \underbrace{\arg \max_{Q \in \mathcal{Q}} H_Q(X)}_{\text{small variance}} + \underbrace{\mathbb{E}_{x \sim Q} [\ln P(x)]}_{\text{assign high/low prob to regions where } P \text{ is large/small}}$$



M-projection

let $P(x_1 \dots x_n)$ be a distribution and \mathcal{Q} be the family of distributions over empty graph G_{\emptyset}

$$\text{then } Q^M = \underset{Q \in \mathcal{Q}}{\operatorname{arg\,min}} D(P||Q) = P(x_1 \dots x_n)$$

$$D(P||Q) = \mathbb{E}_{x \sim P} [\ln P(x) - \ln Q(x)]$$

$$= \mathbb{E}_{x \sim P} [\ln P(x)] - \sum_i \mathbb{E}_{x \sim P} [\ln Q(x_i)]$$

$$= \mathbb{E}_{x \sim P} [\ln \frac{P(x)}{P(x_1 \dots x_n)}] + \sum_i \mathbb{E}_{x \sim P} [\ln \frac{P(x_i)}{Q(x_i)}]$$

$$= D(P||Q^M) + \sum_i D(P(x_i)||Q(x_i))$$

$$\geq D(P||Q^M) \quad \text{equally holds only when } P(x_i) = Q(x_i)$$

$$\begin{aligned} & \int_{\mathbb{R}^n} P(x_1 \dots x_n) \ln Q(x_i) - Q(x_i) dx_1 \dots dx_n \\ &= \sum_i \int_{\mathbb{R}^n} P(x_1 \dots x_n) \ln Q(x_i) dx_1 \dots dx_n \\ &\geq \sum_i \int_{\mathbb{R}} P(x_i) \ln Q(x_i) dx_i \end{aligned}$$

the M-Projection onto factored distribution is simply the product of marginals

Let P be a distribution over x . And \mathcal{Q} be an exponential family defined by $t(\theta)$ and $T(X)$

if there is a parameter θ so $\mathbb{E}_{x \sim \theta} [T(x)] = \mathbb{E}_P [T(X)]$, then M-Projection of P is θ

$$\text{Suppose } \mathbb{E}_P [T(X)] = \mathbb{E}_{x \sim \theta} [T(x)]$$

$$D(P||\theta) - D(P||Q_\theta) = \mathbb{E}_P [\ln \frac{P(x)}{Q_\theta(x)}] - \mathbb{E}_P [\ln \frac{P(x)}{ess(\theta)}]$$

$$= \mathbb{E}_P [\ln Q_\theta(x)] - \mathbb{E}_P [\ln ess(\theta)]$$

$$= \mathbb{E}_P [\ln \frac{1}{Z(t(\theta))} \exp(t(\theta), T(X))] - \mathbb{E}_P [\ln \frac{1}{Z(t(\theta))} \exp(t(\theta), T(X))]$$

$$= \ln(\frac{1}{Z(t(\theta))}) + \langle t(\theta), \mathbb{E}_{x \sim \theta} [T(X)] \rangle - \langle t(\theta), \mathbb{E}_P [T(X)] \rangle$$

$$= \ln(\frac{1}{Z(t(\theta))}) + \langle t(\theta) - t(\theta), \mathbb{E}_{x \sim \theta} [T(X)] \rangle$$

$$= \ln(\frac{1}{Z(t(\theta))}) + \langle t(\theta) - t(\theta), \mathbb{E}_{x \sim \theta} [T(X)] \rangle$$

$$= \int_X Q_\theta(x) \cdot \ln \frac{1}{Z(t(\theta))} \exp(t(\theta), T(X)) dx$$

$$= D(\theta || Q_\theta) \geq 0$$

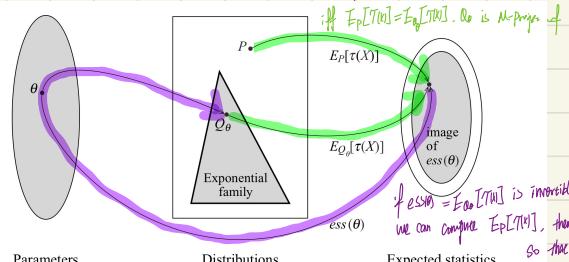
instead of describing a distribution in the family by parameter θ ,

we can describe it in terms of the expected sufficient statistics.

define a mapping from parameter to expected sufficient statistics: $ess(\theta) = \mathbb{E}_{\theta} [T(X)]$

if $\mathbb{E}_P [T(X)]$ is in the image of ess , then the M-projection of P is θ so that $\mathbb{E}_{\theta} [T(X)] = \mathbb{E}_P [T(X)]$

if ess is invertible, parameters of M-projection of P are simply inverses of $ess(\cdot)$



In many exponential families, the sufficient statistics $T(x)$ are moments. (e.g. when $\mathcal{Q} \sim \text{N}(\mu, \sigma^2)$)
when $E_p[T(w)] = E_\theta[T(x)]$. Θ preserves the moments $E[X]$ of P

$$T(w) = \begin{bmatrix} -\frac{w_1}{\sigma^2} \\ w_2 \\ \vdots \\ w_n \end{bmatrix}, T(x) = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$