

1 (**Differential entropy, cross-entropy, KL divergence, and mutual information**) Throughout, use natural logarithms (units: nats), so  $\log(e) = 1$ . Let  $X \in \mathbb{R}^n$  be a continuous random vector with density  $p_X$ . Assume  $\mathbb{E}[X] = \mu$  and  $\text{Cov}(X) = \Sigma \succ 0$ .

(a) Recall the definitions

$$H(X) := - \int p_X(x) \log p_X(x) dx, \quad H(p, q) := - \int p(x) \log q(x) dx,$$

$$\text{KL}(p \| q) := \int p(x) \log \frac{p(x)}{q(x)} dx.$$

Show that

$$H(p, q) = H(p) + \text{KL}(p \| q).$$

(b) Let  $(X, Y)$  have joint density  $p_{X,Y}$  and marginals  $p_X, p_Y$ . Define mutual information by

$$I(X; Y) := \text{KL}(p_{X,Y} \| p_X p_Y).$$

Show the equivalent forms

$$I(X; Y) = H(X) + H(Y) - H(X, Y) = H(X) - H(X | Y) = H(Y) - H(Y | X).$$

(c) Let  $G \sim \mathcal{N}(\mu, \Sigma)$ . Compute  $H(G)$ .

(d) (Gaussian is the maximum-entropy distribution.) Among all random vectors  $X$  with  $\mathbb{E}[X] = \mu$  and  $\text{Cov}(X) = \Sigma$ , show that the Gaussian  $\mathcal{N}(\mu, \Sigma)$  has the largest differential entropy. Equivalently, prove that for any such  $X$ ,

$$H(X) \leq \frac{1}{2} \log((2\pi e)^n \det \Sigma),$$

with equality if and only if  $X \sim \mathcal{N}(\mu, \Sigma)$ .

(e) (Why a constraint is necessary.) Show that differential entropy has no global maximum on  $\mathbb{R}^n$ : if  $X$  has a density and  $a > 0$ , define  $X_a := aX$ . Prove

$$H(X_a) = H(X) + n \log a,$$

and conclude that  $\sup h(X) = +\infty$  if you do not constrain (for example) the covariance.

(f) Let  $P = \mathcal{N}(\mu_0, \Sigma_0)$  and  $Q = \mathcal{N}(\mu_1, \Sigma_1)$  on  $\mathbb{R}^n$  with  $\Sigma_0, \Sigma_1 \succ 0$ .

(f1) Derive a closed form for  $\text{KL}(P \| Q)$ .

(f2) Using (a), write the cross-entropy  $H(P, Q)$  in closed form.

(g) Let  $X \sim \mathcal{N}(0, \Sigma_X)$ ,  $Z \sim \mathcal{N}(0, \Sigma_Z)$  independent, and  $Y := X + Z$ . Compute  $I(X; Y)$  in closed form.

$$(a) H(p) = - \int p(x) \log p(x) dx , KL(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx = \int p(x) (\log p(x) - \log q(x)) dx$$

$$\Rightarrow KL(p||q) + H(p) = \int p(x) \log p(x) - \int p(x) \log q(x) dx = - \int p(x) \log q(x) dx = H(p, q)$$

(b)

$$I(X;Y) = \iint p(x,y) \log \frac{p(x,y)}{p(x)p(y)} dx dy = \iint p(x,y) \log p(x,y) dx dy - \iint p(x,y) \log p(x) dx dy - \iint p(x,y) \log p(y) dx dy$$

$$= -H(X,Y) - \int p(x) \log p(x) dx - \int p(y) \log p(y) dy = -H(X,Y) + H(X) + H(Y)$$

$$H(X,Y) = - \iint p(x,y) \log p(x,y) dx dy = - \iint p(x,y) \log (p(x) \cdot p(y|x)) dx dy = - \iint p(x,y) \log p(x) dx dy - \iint p(x,y) \log p(y|x) dx dy$$

$$= - \int p(x) \log p(x) dx - \iint p(x) p(y|x) \log p(y|x) dx dy = - \int p(x) \log p(x) dx - \int p(x) (\int p(y|x) \log p(y|x) dy) dx = H(X) + H(Y|X) = H(Y) + H(X)$$

$$\Rightarrow I(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

(c)

$$G \sim N(\mu, \Sigma) \Rightarrow p(G) = \frac{1}{(2\pi)^n |\det \Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(G-\mu)^T \Sigma^{-1}(G-\mu)) \Rightarrow \log p(G) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det \Sigma) - \frac{1}{2}(G-\mu)^T \Sigma^{-1}(G-\mu)$$

$$\Rightarrow H(G) = -E \log p(G) = \frac{n}{2} \log(2\pi) + \frac{1}{2} \log(\det \Sigma) + \frac{1}{2} E[(G-\mu)^T \Sigma^{-1}(G-\mu)]$$

$$\text{here } E[(G-\mu)^T \Sigma^{-1}(G-\mu)] = E[\text{tr}((G-\mu)^T \Sigma^{-1}(G-\mu))] = E[\text{tr}(\Sigma^{-1}(G-\mu)(G-\mu)^T)] = \text{tr}(\Sigma^{-1} E[(G-\mu)(G-\mu)^T]) = \text{tr}(\Sigma^{-1} \Sigma) = n$$

$$\Rightarrow H(G) = \frac{n}{2} \log(2\pi) + \frac{1}{2} \log(\det \Sigma) + \frac{n}{2} = \frac{1}{2} \log((2\pi e)^n \det \Sigma)$$

(d)

$$E[X] = \mu, \text{Cov}(X) = \Sigma, \quad KL(p_x||p_a) = -H(X) + H(p_x, p_a) \geq 0, \quad H(p_x, p_a) = - \int p_x(x) \log p_a(x) dx$$

$$\text{since } \log p_a(x) = C - \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu) \Rightarrow \text{only depend on } \mu \text{ and } \Sigma \text{ of } X, \text{ so } \int p_x(x) \log p_a(x) dx = \int p_a(x) \log p_a(x) dx$$

$$\Rightarrow H(p_x, p_a) = H(G) \Rightarrow H(X) \leq H(G) = \frac{1}{2} \log((2\pi e)^n \det \Sigma)$$

(e)

$$\text{the density of } X_a \text{ is } p_{X_a}(y) = \frac{1}{a^n} p_X(\frac{y}{a}) \Rightarrow H(X_a) = - \int p_{X_a}(y) \log p_{X_a}(y) dy = - \int \frac{1}{a^n} p_X(\frac{y}{a}) \log(\frac{1}{a^n} p_X(\frac{y}{a})) dy$$

$$\text{here } x = \frac{y}{a}, dy/a^n dx \Rightarrow H(X_a) = - \int p_X(x) (\log p(x) - \log a^n) dx = - \int p_X(x) \log p(x) dx + n \log a \int p(x) dx = H(X) + n \log a$$

then if  $a \rightarrow \infty$ ,  $H(X_a) \rightarrow +\infty$ ,  $\text{Cov}(X_a) = a^2 \text{Cov}(X) = a^2 \Sigma$ , so  $\sup h(x) = +\infty$  if covariance not constrained

(f)

$$(f) KL(P||Q) = E_p[\log P - \log Q] = E_p[\log P] - E_p[\log Q]$$

$$E_p[\log P] = -H(P) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det \Sigma_0) - \frac{n}{2}$$

$$\log Q = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det \Sigma_1) - \frac{1}{2}(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1)$$

$$\text{here } E_p[(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1)] = E_p[\text{tr}(\Sigma_1^{-1}(x-\mu_1)(x-\mu_1)^T)] = \text{tr}(\Sigma_1^{-1} E_p[(x-\mu_0 + \mu_0 - \mu_1)(x-\mu_0 + \mu_0 - \mu_1)^T]), \text{ here } E_p[x-\mu_0] = 0$$

$$\text{so } E_p[(x-\mu_1)^T \Sigma_1^{-1}(x-\mu_1)] = \text{tr}(\Sigma_1^{-1} [\Sigma_0 + (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T]) = \text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_0 - \mu_1)^T \Sigma_1^{-1} (\mu_0 - \mu_1)$$

$$\Rightarrow E_p[\log Q] = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log(\det \Sigma_1) - \frac{1}{2} [\text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_0 - \mu_1)^T \Sigma_1^{-1} (\mu_0 - \mu_1)]$$

$$\Rightarrow KL(P||Q) = E_p[\log P] - E_p[\log Q] = \frac{1}{2} (\log \frac{\det \Sigma_1}{\det \Sigma_0} - n + \text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_0 - \mu_1)^T \Sigma_1^{-1} (\mu_0 - \mu_1))$$

$$(f2) H(P, Q) = H(P) + KL(P||Q) = \frac{n}{2} \log(2\pi) + \frac{1}{2} \log(\det \Sigma_0) + \frac{n}{2} + \frac{1}{2} (\log \frac{\det \Sigma_1}{\det \Sigma_0} - n + \text{tr}(\Sigma_1^{-1} \Sigma_0) + (\mu_0 - \mu_1)^T \Sigma_1^{-1} (\mu_0 - \mu_1))$$

$$= \frac{1}{2} (\ln \log(2\pi) + \log(\det \Sigma_1) + \text{tr}(\Sigma_1^{-1} \Sigma_0) + \text{tr}(\Sigma_2^{-1} \Sigma_0) + (\mu_0 - \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1))$$

$$(g) I(X; Y) = H(Y) - H(Y|X), \quad Y|X=x \sim (X+2)|X=x \quad \Rightarrow \quad H(Y|X) = H(2) = \frac{1}{2} \log((2\pi e)^n \det \Sigma_2)$$

$$Y \sim N(0, \Sigma_X + \Sigma_Z) \Rightarrow H(Y) = \frac{1}{2} \log((2\pi e)^n \det(\Sigma_X + \Sigma_Z))$$

$$\Rightarrow I(X; Y) = \frac{1}{2} \log((2\pi e)^n \det(\Sigma_X + \Sigma_Z)) - \frac{1}{2} \log((2\pi e)^n \det \Sigma_2) = \frac{1}{2} \log \frac{\det(\Sigma_X + \Sigma_Z)}{\det \Sigma_2} = \frac{1}{2} \log \det(I_n + \Sigma_2^{-1} \Sigma_X)$$