

## ■ Common formulae of information theory

Hugo Touchette (htouchet@mit.edu)

Updated: August 31, 2004

Taken from T. Cover, J. Thomas, *Elements of Information Theory*, Wiley, 1991.

**Entropy.**  $H(X) = - \sum_x p(x) \log p(x)$

$$H(X) \geq 0$$

$$H_b(X) = \log_b a H_a(X)$$

$$\log_a x = \log_a b \log_b x = \frac{\log_b x}{\log_b a}$$

$$H(X) \leq \log |\mathcal{X}|$$

$$d_\alpha \sum_x p(x)^\alpha \Big|_{\alpha=1} = H(X)$$

$$(a^x)' = a^x \ln a, \quad (\log_a x)' = \frac{1}{x \ln a} = \frac{\log_a e}{x}$$

**Joint entropy.**  $H(X, Y) = - \sum_{x,y} p(x, y) \log p(x, y)$

$$H(X_1, \dots, X_n) \leq \sum_{i=1}^n H(X_i) \text{ eq. iff } X_i \text{ ind.}$$

**Conditional entropy.**  $H(Y|X) = \sum_x H(Y|x)p(x) = - \sum_{x,y} p(x, y) \log p(y, x)$

$$H(Y|X) \geq 0 \text{ eq. if } Y = f(X), \text{ or iff } Y|x \text{ is deterministic for all } x \in \text{supp}(X)$$

$$H(Y|X) \leq H(Y) \text{ eq. iff } X, Y \text{ ind.}$$

**Relative entropy.**  $D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$

$$D(p||q) \geq 0$$

$$D(p||u) = \log |\mathcal{X}| - H(X), \quad u(x) = |\mathcal{X}|^{-1}$$

**Mutual information.**  $I(X; Y) = \sum_{x,y} p(x, y) \log \frac{p(x, y)}{p(x)p(y)} = D(p(x, y)||p(x)p(y))$

$$I(X; Y) \geq 0 \text{ eq. iff } X, Y \text{ ind.}$$

$$I(X; Y) = I(Y; X)$$

$$I(X; X) = H(X)$$

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

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

$$\max I(X; Y) = \min(H(X), H(Y))$$

**Information metric.**  $\Delta(X, Y) = H(X|Y) + H(Y|X)$

$$\Delta(X, Y) = H(X) + H(Y) - 2I(X; Y) = H(X, Y) - I(X, Y)$$

**Chain rules.**

$$H(X, Y) = H(X) + H(Y|X)$$

$$H(X, Y|Z) = H(X|Z) + H(Y|X, Z)$$

$$H(X_1, \dots, X_n) = \sum_{i=1}^n H(X_i|X_{i-1}, \dots, X_1)$$

$$I(X_1, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y | X_{i-1}, X_{i-2}, \dots, X_1)$$

$$D(p(x, y) || q(x, y)) = D(p(x) || q(x)) + D(p(y|x) || q(y|x))$$

### Conditioning.

$$D(p(y|x) || q(y|x)) = \sum_x p(x) \sum_y p(y|x) \log \frac{p(y|x)}{q(y|x)} \quad \text{Conditional relative entropy}$$

$$I(X; Y | Z) = H(X|Z) - H(X|Y, Z) \quad \text{Conditional mutual information}$$

### Convexity properties.

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2) \text{ or } f''(x) \geq 0 \quad \text{Convex function}$$

$$E[f(X)] \geq f(E[X]) \quad \text{Jensen's inequality}$$

$$\sum_{i=1}^n a_i \log \frac{a_i}{b_i} \geq \left( \sum_{i=1}^n a_i \right) \log \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i} \quad \text{log sum inequality}$$

### Differential entropy.

$$h(X) = - \int p(x) \log p(x) dx$$

$$h(X + a) = h(X)$$

$$H(X^\Delta) + \log \Delta \rightarrow h(X)$$

$$h(AX) = h(X) + \log |A|$$

$$D(P || Q) \geq 0$$

$$h(X) \leq \log \text{supp } X$$

$$D(P || Q) = \sup_{\Delta} D(P^\Delta || Q^\Delta)$$

$$h(X) = \frac{1}{2} \log 2\pi e \sigma^2 \text{ when } X \sim N(m, \sigma^2)$$

### Correlations and causation.

$$A \rightarrow B \leftarrow C \quad I(A; C) = 0$$

$$I(A; C | B) \neq 0 \text{ in general}$$

$$A \leftarrow B \rightarrow C \text{ or } I(A; C) \neq 0 \text{ in general}$$

$$A \rightarrow B \rightarrow C \quad I(A; C | B) = 0 \text{ (bottleneck)}$$

### Chains of random variables.

$$X \rightarrow Y \rightarrow Z \quad I(X; Y) \geq I(X; Z)$$

$$I(Z; Y) \geq I(Z; X) \text{ eq. iff } I(X; Y | Z) = 0$$

$$I(X; Y) \geq I(X; g(Y))$$

$$I(X; Y | Z) \leq I(X; Y)$$

$$X^{(n)} \rightarrow Y^{(k)} \rightarrow Z^{(m)} \quad I(X; Z) \leq \log k$$

$$k < n, k < m \quad I(X; Z) = 0 \text{ if } k = 1$$

### Asymptotic equipartition theorem.

$$X_1 X_2 \dots X_n \sim p(x) \text{ iid.}$$

$$-\frac{1}{n} \log p(x_1, x_2, \dots, x_n) = -\frac{1}{n} \sum_{i=1}^n \log p(x_i)$$

$$\rightarrow -E[\log p(x_i)] = H(X) \text{ (in probability)}$$

$$A_\varepsilon^n = \{x^n \in \mathcal{X}^n : 2^{-n(H+\varepsilon)} \leq p(x^n) \leq 2^{-n(H-\varepsilon)}\} \quad \text{Typical set}$$

$$-\frac{1}{n} \log p(x^n) = H(X) \text{ (within } \varepsilon)$$

$$\Pr\{A_\varepsilon^n\} > 1 - \varepsilon$$

$$\text{(from above result)}$$

$$|A_\varepsilon^n| \leq 2^{n(H+\varepsilon)}$$

$$|A_\varepsilon^n| \doteq 2^{nH} \text{ (within } \varepsilon \text{ exponentially)}$$

$$|A_\varepsilon^n| \geq (1 - \varepsilon) 2^{n(H-\varepsilon)}$$

$$p(x^n) \doteq 2^{-nH}$$

### Method of types.

|   |  |                  |
|---|--|------------------|
| $X_1 X_2 \dots X_n, x_1 x_2 \dots x_n = x^n = \mathbf{x} \in \mathcal{X}^n$                                 | $\mathcal{X} = \{a_1, a_2, \dots, a_{ \mathcal{X} }\}$ |                  |
| $P_{\mathbf{x}}(a) = N(a \mathbf{x})/n$   | $\sum_a P_{\mathbf{x}}(a) = 1$                         | Type             |
| $\mathcal{P}_n = \{P_{\mathbf{x}} :  \mathbf{x}  = n\}$   |  | Set of types $n$ |
| $T(P) = \{\mathbf{x} \in X^n : P_{\mathbf{x}} = P\}$  |  | Type class       |
| $ \mathcal{P}_n  \leq (n+1)^{ \mathcal{X} }$  | $ \mathcal{X}^n  \sim  \mathcal{X} ^n$                 |                  |
|   | $ \mathcal{P}_n  \sim n^{ \mathcal{X} }$               |                  |
| $Q^n(\mathbf{x}) = 2^{-n[H(P_{\mathbf{x}}) + D(P_{\mathbf{x}}  Q)]}$  | $X_1 X_2 \dots X_n \sim \text{iid } Q(x)$              |                  |
| $Q^n(\mathbf{x}) = 2^{-nH(Q)}$  | $\mathbf{x} \in T(Q)$                                  | Type $Q$ seqs.   |
| $\frac{1}{(n+1)^{ \mathcal{X} }} 2^{nH(P)} \leq  T(P)  \leq 2^{nH(P)}$                                      | $ T(P)  \doteq 2^{nH(P)}$                              | Type class size  |
| $\frac{1}{(n+1)^{ \mathcal{X} }} 2^{-nD(P  Q)} \leq Q^n(T(P)) \leq 2^{-nD(P  Q)}$                           | $Q^n(T(P)) \doteq 2^{-nD(P  Q)}$                       | Type class prob. |
| $\Pr\{D(P_{\mathbf{x}}  Q) > \varepsilon\} \leq 2^{-n[\varepsilon -  \mathcal{X}  \frac{\log n + 1}{n}]}$   | $X_1 X_2 \dots X_n \sim \text{iid } Q(x)$              | WLLN             |
| $\Pr\{D(P_{\mathbf{x}}  Q) > \varepsilon\} \leq n^{ \mathcal{X} } 2^{-n\varepsilon} \sim 2^{-n\varepsilon}$ | $D(P_{\mathbf{x}}  Q) \rightarrow 0$ (ip)              |                  |
| $Q^n(E) \leq (n+1)^{ \mathcal{X} } 2^{-nD(P^*  Q)}$   | $X_1 X_2 \dots X_n \sim \text{iid } Q(x)$              | Sanov            |
| $P^* = \arg \min_{P \in E} D(P  Q)$   | $P \subseteq \mathcal{P}$                              |                  |

### Rate distortion theory.

|   |  |                      |
|---|--|----------------------|
|   | $X^n \rightarrow f(X^n) \rightarrow g(f(X^n)) \rightarrow \hat{X}^n$ |                      |
| $d : \mathcal{X} \times \hat{\mathcal{X}} \rightarrow \mathbb{R}^+$   | $d(x^n, \hat{x}^n) = \frac{1}{n} \sum_i d(x_i, \hat{x}_i)$           | Distance, distortion |
|   | $\max d(x, \hat{x}) < \infty$  |                      |
| $d = \begin{cases} 0, & x=\hat{x} \\ 1, & x \neq \hat{x} \end{cases}$   | $E[d(X, \hat{X})] = \Pr\{X \neq \hat{X}\}$                           | Hamming distance     |
| $R(D) = \min_{p(x x): E[d(X^n, \hat{X}^n)] \leq D} I(X; \hat{X})$   | $f(X^n) \in \{1, 2, \dots, 2^{nR}\}$                                 | Rate (# bits needed) |
| $R(D) = \begin{cases} \frac{1}{2} \log \frac{\sigma^2}{D}, & 0 \leq D \leq \sigma^2 \\ 0, & D > \sigma^2 \end{cases}$ |  | Gaussian channel     |

### Elements of probability theory.

|  |   |
|--|---|
| $X_i \sim \text{iid}, \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\text{ip}} E[X]$ | WLLN  |
| $Y = g(X), X = h(Y) = g^{-1}(Y)$   | $f_Y(y) = f_X(h(y))  h'(y)  = f_X(x) \left  \frac{\partial h(y)}{\partial y} \right $         |
| $Y = g(X) = \alpha X + \beta$  | $\mathcal{N}(\mu, \sigma^2) \xrightarrow{g} \mathcal{N}(\alpha\mu + \beta, \alpha^2\sigma^2)$ |
| $Y = X + Z$  | $X \sim \mathcal{N}(x, \sigma_X^2)$   |
|  | $Z \sim \mathcal{N}(z, \sigma_Z^2)$   |
|  | $Y \sim \mathcal{N}(x + z, \sigma_X^2 + \sigma_Z^2)$  |