

## §1. Entropy, Divergence, and Mutual Information

$$H(X) := -\sum_x p_X(x) \log p_X(x) = -\mathbb{E}(\log p_X(x)) \quad D(p||q) := \sum_x p(x) \frac{p(x)}{q(x)} = \mathbb{E}(\log \frac{1}{q(X)}) - H(X)$$

$$I(X;Y) := \sum_{x,y} p_{X,Y}(x,y) \log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)} = D(p_{X,Y}||p_X p_Y) = H(X) + H(Y) - H(X,Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

$$H(Y|X) = -\sum_{x,y} p_{X,Y}(x,y) \log p_{Y|X=x}(y) = -\sum_x p_X(x) \mathbb{E}(\log p_{Y|X=x}(Y)) = -\mathbb{E}(\log p_{Y|X}(Y)) = \sum_x H(Y|X=x) \mathbb{P}(X=x) = H(X,Y) - H(X)$$

$$I(X;Y|Z) := H(X|Z) - H(X|Y,Z)$$

$$\text{Gibb's Inequality } H(X) = -\sum_x p(x) \log p(x) \leq -\sum_x p(x) \log q(x) \quad (= \text{iff } p = q)$$

$$D(p_{X,Y}||p_{\hat{X},\hat{Y}}) = D(p_{Y|X}||p_{\hat{Y}|\hat{X}}|p_X) + D(p_X||p_{\hat{X}})$$

$$D(p_{Y|X}||p_{\hat{X},\hat{Y}}|p_X) = D(p_X p_{Y|X}||p_X p_{\hat{Y}|\hat{X}})$$

$$D(p_{Y|X}||q_{Y|X}|p_X) = \sum_x p_X(x) D(p_{Y|X=x}||q_{Y|X=x})$$

$$\text{Logsum: } a_i, b_i \geq 0 \implies \sum_{i=1}^n a_i \log \left( \frac{a_i}{b_i} \right) \geq \left( \sum_{i=1}^n \log \left( \frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i} \right) \right) \text{ Pf } q_i = a_i / \sum_{i=1}^n a_i \text{ and Gibbs}$$

$$D(\lambda p_1 + (1-\lambda)p_2||\lambda q_1 + (1-\lambda)q_2) \leq \lambda D(p_1||q_1) + (1-\lambda) D(p_2||q_2) \text{ for } \lambda \in [0,1] \text{ (Pf: logsum)}$$

$$I(X;Y) \geq 0 \quad (= \text{iff } X \perp Y) \quad I(X;Y) = H(X) - H(X|Y)$$

$$H(X,Y) = H(X) + H(Y|X) \quad I(X_1, \dots, X_n; Y) = \sum_{i=1}^n I(X_i; Y|X_{i-1}, \dots, X_1) \text{ (Chain)}$$

$$\text{Data Processing } (X \perp Z)|Y \implies I(X;Y) \geq I(X;Z) \text{ Pf Chain } I(Y,Z;X) \text{ and } I(Z,Y;X)$$

$$f: \mathcal{X} \rightarrow \mathcal{Y} \implies I(X;Y) \geq I(X;f(Y))$$

$$0 \leq H(X) \leq \log(|\mathcal{X}|) \quad (= \text{iff } X \text{ const, } X \text{ uniform respectively})$$

$$0 \leq H(X|Y) \leq H(X) \quad (2^{\text{nd}} = \text{iff } X \perp Y \text{ iff } X = f(Y))$$

$$H(f(X)) \leq H(X) \quad (= \text{iff } f \text{ bijective})$$

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

$$X, Y \text{ iid} \implies \mathbb{P}(X=Y) \geq 2^{-H(X)} \quad (= \text{iff uniform, proof by Jensen})$$

$$\text{Fano's Inequality } H(X|Y) \leq H(\mathbf{1}_{X \neq Y}) + \mathbb{P}(X \neq Y) \log(|\mathcal{X}| - 1) \quad (\text{Proof by defining } Z = \mathbf{1}_{X \neq Y} \text{ and } H(Z|X,Y) = 0, \text{ so that } H(X|Y) = H(Z|Y) + H(X|Y,Z))$$

$$\text{Cor: } H(X|Y) \leq 1 + \mathbb{P}(X \neq Y) \log(|\mathcal{X}| - 1)$$

Max Entropy Tip: Use Gibb's to bound  $H(q)$ , and equality iff  $p = q$  where  $p$  specified as:

$$\text{For } \mathbb{E}(X) = \text{const. on } \{1, \dots\}, X \sim \text{geom}(p) \in \{1, \dots\}, \mathbb{E}(X) = \frac{1}{p}, H(X) = \frac{-(1-p) \log 1-p - p \log p}{p}$$

$$\text{For } \mathbb{E}(X) = \text{const. on } \{0, \dots\}, X \sim \text{geom}(p) \in \{0, \dots\}, \mathbb{E}(X) = \frac{1-p}{p}, H(X) = \frac{-(1-p) \log 1-p - p \log p}{p}$$

$$\text{Max Entropy for } \mathbb{E}(f(X)) = C, X \sim p_X(x) = \frac{e^{-\lambda f(x)}}{\sum_x e^{-\lambda f(x)}} \text{ where } \lambda \text{ chosen s.t. } \mathbb{E}(f(X)) = C$$

$$\text{§2. AEP WLLN: } \bar{X}_n \xrightarrow{P} \mu \text{ (ie: } \lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| < \epsilon) = 1 \text{ for any } \epsilon > 0)$$

$$\text{Weak AEP 1: } -\frac{1}{n} \log p_{X_1, \dots, X_n}(X_1, \dots, X_n) \xrightarrow{P} H(X)$$

$$\mathcal{T}_n^\epsilon := \{(x_1, \dots, x_n) \in \mathcal{X}^n : |-\frac{1}{n} \log p_{X_1, \dots, X_n}(x_1, \dots, x_n) - H(X)| \leq \epsilon\}$$

$$\forall \epsilon > 0, \exists N \text{ s.t. } \forall n \geq N,$$

$$(i) p_{X_1, \dots, X_n}(x_1, \dots, x_n) \in [2^{-n(H(X)+\epsilon)}, 2^{-n(H(X)-\epsilon)}] \text{ for any } (x_1, \dots, x_n) \in \mathcal{T}_n^\epsilon$$

$$(ii) \mathbb{P}((X_1, \dots, X_n) \in \mathcal{T}_n^\epsilon) \geq 1 - \epsilon \text{ (From weak AEP 1)}$$

$$(iii) |\mathcal{T}_n^\epsilon| \in [(1-\epsilon)2^{n(H(X)-\epsilon)}, 2^{n(H(X)+\epsilon)}] \quad \text{Pf } 1 = \sum_x p_{\mathbf{X}}(\mathbf{x}) \geq \sum_{\mathcal{T}_n^\epsilon} p_{\mathbf{X}}(\mathbf{x}) \geq \sum_{\mathcal{T}_n^\epsilon} 2^{-n(H(X)+\epsilon)} \text{ and } 1 - \epsilon \leq \mathbb{P}((X_1, \dots, X_n) \in \mathcal{T}_n^\epsilon) \leq \sum_{\mathcal{T}_n^\epsilon} 2^{-n(H(X)-\epsilon)} = 2^{-n(H(X)-\epsilon)} |\mathcal{T}_n^\epsilon|$$

$$\text{Shannon's 1st Thrm: } \forall \epsilon > 0, \exists n \in \mathbb{Z} \text{ and } c: \mathcal{X}^* \rightarrow \{0,1\}^* \text{ s.t. } \cup_{k \geq 0} \mathcal{X}^{nk} \rightarrow \{0,1\}^*, (x_1, \dots, x_k) \rightarrow c(x_1) \dots c(x_k) \in \{0,1\}^* \text{ injective and } \frac{1}{n} \mathbb{E}(|c(X_1, \dots, X_n)|) \leq H(X) + \epsilon$$

## §3. Optimal Codes

$$\lceil A \rceil - 1 < A \leq \lceil A \rceil \quad A \leq \lceil A \rceil < A + 1$$

$$\text{Kraft-McMillan: (i) } c: \mathcal{X} \rightarrow \mathcal{Y}^* \text{ uniquely decodable and } l_x := |c(x)| \rightarrow \sum_{x \in \mathcal{X}} |\mathcal{Y}|^{-l_x} \leq 1$$

$$\text{Pf) } (\sum_x d^{-|c(x)|})^n = \sum_{k=nl_{\min}}^{nl_{\max}} a(k) d^{-k}, a(k) < d^k \text{ by uniq decod. Take } n^{\text{th}} \text{ root and } n \rightarrow \infty.$$

$$(ii) \text{ Given } (l_x)_{x \in \mathcal{X}} \subset \mathbb{N} \text{ and } \sum_{x \in \mathcal{X}} |\mathcal{Y}|^{-l_x} \leq 1, \text{ then } \exists \text{ prefix code } c: \mathcal{X} \rightarrow \mathcal{Y}^* \text{ s.t. } |c(x)| = l_x$$

$$\text{Pf) Relabel } \mathcal{X} = \{1, \dots, |\mathcal{X}|\}, l_1 \leq \dots \leq l_{|\mathcal{X}|}. r_m := \sum_{i=1}^{m-1} |\mathcal{Y}|^{-l_i}. c(m): \text{ first } m \text{ digits of } r_m$$

$$X \text{ be RV in finite } \mathcal{X} \text{ and } c \text{ (uniq decod., } d\text{-ary), then } H_d(X) \leq \mathbb{E}(|c(X)|) (= \text{iff } |c(x)| = -\log_x p_X(x),$$

$$\text{Lower bd on length) Pf) Let } l_x = c(x), q(x) = \frac{d^{-l_x}}{\sum_{x \in \mathcal{X}} d^{-l_x}}, \text{ consider } \mathbb{E}(|c(X)|) - H_d(X)$$

$$\text{Existence of Optimal Code } H_d(X) \leq \mathbb{E}(|c^*(X)|) < H_d(X) + 1 \quad \text{Pf) } l_x = \lceil -\log_d(p(x)) \rceil \text{ Kraft-McMillan}$$

1	<b>Shannon's Code</b> (i) Order $p_1 \geq \dots \geq p_m$ . (ii) $c(x_r) = \text{first } l_r := \lceil -\log_{ \mathcal{Y} }(p_r) \rceil$ digits of $\sum_{i=1}^{r-1} p_i$
2	<u>Shannon with distrib. estimation</u> $p, q$ (estimation) pmf on $\mathcal{X}$ ,
3	$H_d(X) + D_d(p  q) \leq \mathbb{E}( c_q(X) ) < H_d(X) + D_d(p  q) + 1$
4	Pf: Bound $\mathbb{E}( c_q(X) ) = \sum_x p(x) \lceil -\log_d(q(x)) \rceil$
5	<b>Elias' code</b> First $\lceil -\log_d(p_i) \rceil + 1$ digits of $\sum_{i < r} p_i + \frac{p_r}{2} \implies H_d(X) + 1 \leq \mathbb{E}( c_E(X) ) \leq H_d(X) + 2$
6	
7	Bijection between $d$ -ary prefix codes and $d$ -ary rooted trees.
8	<b>Huffman code is optimal</b> Pf) if $p_1 \geq \dots \geq p_m$ , (i) $p_j > p_k \implies  c(x_j)  \leq  c(x_k) $ , (ii) two longest
9	codewords have same len (iii) two longest codewords only differ in last digit. $p, p' \implies L(c^p) - L(c^{p'}) =$
10	$p_{m-1} + p_m$ Also for $e^p, e^{p'}$ , $L(e^p) - L(e^{p'}) = p_{m-1} + p_m$ . Subtract each other, $L(e^p) = L(c^p)$
11	
12	<b>§4 Channel Coding</b> DMC $(\mathcal{X}, M, \mathcal{Y})$ where $\mathcal{X}$ (input alphabet), $\mathcal{Y}$ (output alphabet), $M$ ( $ \mathcal{X}  \times  \mathcal{Y} $
13	stochastic matrix).
14	Channel Capacity: $C := \sup I(x; y) = H(Y) - H(Y X)$
15	Tip: $H(Y X) = \sum_x H(Y X=x)p_X(x)$ , Use $I(X, Y) = H(Y) - H(Y X)$
16	$(m, n)$ -channel code for DMC $(\mathcal{X}, M, \mathcal{Y})$ : tuple $(c, d)$ where $c : \{1, \dots, m\} \rightarrow \mathcal{X}^n$ (Encoder) and
17	$d : \mathcal{Y}^n \rightarrow \{1, \dots, m\}$ (Decoder)
18	Rate of $(m, n)$ -code $(c, d)$ : $\rho(c, d) := \frac{1}{n} \log_{ \mathcal{X} }(m)$
19	$\epsilon_i = \mathbb{P}(d(\mathbf{Y} \neq i c(i)) = \mathbf{X})$ for $i = 1, \dots, m$ , $\epsilon_{\max} := \max_{i \in \{1, \dots, m\}} \epsilon_i$ , $\bar{\epsilon} := \frac{1}{m} \sum_{i=1}^m \epsilon_i$
20	Rate $R$ achievable if $\forall \epsilon > 0$ , $\exists$ suff. large. $m, n$ and $(m, n)$ -channel code $(c, d)$ with $\rho(c, d) > R - \epsilon$
21	and $\epsilon_{\max} < \epsilon$
22	<u>Shannon's 2<sup>nd</sup> theorem</u> DMC $(\mathcal{X}, M, \mathcal{Y})$ with capacity $C$ , then $R > 0$ achievable iff $R \leq C$
23	Pf)
24	$\mathcal{J}_\epsilon^{(n)} =$
25	$\left\{ (x, y) \in \mathcal{X}^n \times \mathcal{Y}^n : \max \left( \left  \frac{-\log p_{\mathbf{X}, \mathbf{Y}}(x, y)}{n} - H(X, Y) \right , \left  \frac{-\log(p_{\mathbf{X}}(x))}{n} - H(X) \right , \left  \frac{-\log(p_{\mathbf{Y}}(y))}{n} - H(Y) \right  \right) \right\}$
26	<u>Joint AEP</u> : $\mathbf{X} = (X_1, \dots, X_n)$ , $\mathbf{Y} = (Y_1, \dots, Y_n)$
27	(i) $\lim_{n \rightarrow \infty} \mathbb{P}((\mathbf{X}, \mathbf{Y}) \in \mathcal{J}_\epsilon^{(n)}) = 1$ (ii) $ \mathcal{J}_\epsilon^{(n)}  \leq 2^{n(H(X, Y) + \epsilon)}$
28	(iii) $\exists n_0$ s.t. $\forall n \geq n_0$ , $(1 - \epsilon)2^{-n(I(X; Y) + 3\epsilon)} \leq \mathbb{P}((\mathbf{X}', \mathbf{Y}') \in \mathcal{J}_\epsilon^{(n)}) \leq 2^{-n(I(X; Y) - 3\epsilon)}$
29	
30	<b>Channel Coding w/ non-iid Input</b>
31	<u>Stationary stochastic process</u> : $\mathbb{P}(X_1 = x_1, \dots, X_n = x_n) = \mathbb{P}(X_{1+j} = x_1, \dots, X_{n+j} = x_n)$ for all
32	$n, j \in \mathbb{Z}$
33	<u>Entropy rate of stochastic process</u> : $\mathcal{H}(X) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n)$
34	Lemma (1): Stationary Stochastic process $X$ , $n \rightarrow H(X_n X_{n-1}, \dots, X_1)$ is non-increasing and
35	$\lim_{n \rightarrow \infty} H(X_n X_{n-1}, \dots, X_1)$ exists.
36	Lemma (2): $\lim_{n \rightarrow \infty} a_n = a \implies \frac{1}{n} \sum_{i=1}^n a_i = a$
37	Thrm: Stationary stoch. process $X \implies \mathcal{H}(X) = \lim_{n \rightarrow \infty} H(X_n X_{n-1}, \dots, X_1)$ PF) Lemma (i),(ii)
38	Lemma (3): $H(Y_n Y_{n-1}, \dots, Y_1)$
39	Lemma (4): $H(Y_n Y_{n-1}, \dots, Y_2, Y_1) \leq \lim_k H(Y_{n+k+1} Y_{n+k}, \dots, Y_1) = \mathcal{H}(Y)$
40	Thrm: $X(X_i)_{i \geq 1}$ stationary Markov. $\phi : \mathcal{X} \rightarrow \mathcal{Y}$ , $Y_i := \phi(X_i)$ , then $H(Y_n Y_{n-1}, \dots, Y_1, X_1) \leq$
41	$\mathcal{H}(Y) \leq H(Y_n Y_{n-1}, \dots, Y_1)$ , $\mathcal{H}(Y) = \lim_{n \rightarrow \infty} H(Y_n Y_{n-1}, \dots, Y_1, X_1) = \lim_{n \rightarrow \infty} H(Y_n Y_{n-1}, \dots, Y_1)$
42	PF) Lemma (3),(4)
43	
44	<b>§Appendix</b>
45	<u>Markov Inequality</u> : $\mathbb{P}(X \geq x) \leq \frac{\mathbb{E}(X)}{x}$ for all $x > 0$
46	<u>Chebyshev Inequality</u> : $\mathbb{P}( X - \mu  > \epsilon) \leq \frac{\sigma^2}{\epsilon^2}$
47	<u>Optimal <math>\neq</math> Huffman</u> : $0.3 = 00, 0.3 = 10, 0.2 = 01, 0.2 = 11$