

15 Entropies and hypothesis testing

15.1 Entropies

15.1.1 Notation

Given more than one random variable, X and Y for example,

$$H(X, Y) := S(P_{XY})$$

denotes the “joint entropy” of X and Y . Note that this is simply the entropy of the random variable (X, Y) .

If we have states ρ_{ABC} and σ_{ABD} then, unless otherwise specified, $\rho_A := \text{Tr}_{BC}\rho_{ABC}$, $\rho_{AC} := \text{Tr}_B\rho_{ABC}$, $\sigma_{BD} := \text{Tr}_A\sigma_{ABD}$, etc. denote the states of the various subsystems. In analogy with our notation for entropies of random variables we write the von Neumann entropies for the states of these systems $H(A)_\rho := S(\rho_A)$, $H(D)_\sigma := S(\sigma_D)$, $H(ABC)_\rho := S(\rho_{ABC})$, $H(AB)_\rho := S(\rho_{AB})$, $H(AB)_\sigma := S(\sigma_{AB})$, and so on.

15.1.2 Classical information stored in quantum systems

Suppose we have a number of random variables, e.g. X, Y and Z , whose values are stored in quantum systems \tilde{X}, \tilde{Y} and \tilde{Z} in the standard way, so that the state of $\tilde{X}\tilde{Y}\tilde{Z}$ is

$$\rho_{\tilde{X}\tilde{Y}\tilde{Z}} = \sum_{x,y,z} P_{XYZ}(x, y, z) |x\rangle\langle x|_{\tilde{X}} \otimes |y\rangle\langle y|_{\tilde{Y}} \otimes |z\rangle\langle z|_{\tilde{Z}}. \quad (15.1)$$

The quantum entropies of these systems correspond to the classical entropies of the RVs in the obvious way e.g. $H(\tilde{X}\tilde{Y}\tilde{Z})_\rho = H(X, Y, Z)$, $H(\tilde{Y}\tilde{Z})_\rho = H(Y, Z)$, $H(\tilde{X})_\rho = H(X)$ etc.

15.1.3 Basic bounds on entropy

Given any distribution p on a finite set \mathcal{A} ,

$$0 \leq S(p) \leq \log |\text{supp}(p)| \leq \log |\mathcal{A}|, \quad (15.2)$$

where the first equality holds iff $p(x) = 1$ for some $x \in \mathcal{A}$. This result can be established by using the strict concavity of \log and Jensen’s inequality. Doing so is an exercise on example sheet 3.

Given any state ρ on a finite dimensional Hilbert space \mathcal{H} ,

$$0 \leq S(\rho) \leq \log \dim(\text{supp}(\rho)) \leq \log \dim(\mathcal{H}), \quad (15.3)$$

where the first equality holds iff ρ is pure. These inequalities follow easily from those in (15.2) by using the fact that if $\rho = \sum_k p(k) |\alpha_k\rangle\langle \alpha_k|$ is an eigendecomposition for ρ then $S(\rho) = S(p)$.

15.2 Conditional entropy and the chain rule

Definition 1 (Classical conditional entropies). Given random variables X and Y and Z :

1. $H(X|Y = y) := S(P_{X|Y=y})$;
2. $H(X|Y) := \sum_{y \in \mathcal{A}_Y} H(X|Y = y) \Pr(Y = y)$;
3. $H(X|Y, Z = z) = \sum_{y \in \mathcal{A}_Y} H(X|Y = y, Z = z) \Pr(Y = y|Z = z)$.

It follows from the results in section 15.1.3 that $H(X|Y) \geq 0$ with equality iff $X = f(Y)$ for some function $f : \text{supp}(P_Y) \rightarrow \mathcal{A}_X$. Similarly, $H(X|Y, Z = z) \geq 0$ with equality iff, given $Z = z$, $X = f(Y)$ for some function $f : \text{supp}(P_{Y|Z=z}) \rightarrow \mathcal{A}_X$.

The product rule of probability says that $P_{XY}(x, y) = P_{X|Y=y}(x)P_Y(y)$, so

$$\log \frac{1}{P_{XY}(x, y)} = \log \frac{1}{P_{X|Y=y}(x)} + \log \frac{1}{P_Y(y)}. \quad (15.4)$$

Multiplying by $P_{XY}(x, y)$ and summing over x and y

$$\begin{aligned} \sum_{x, y} P_{XY}(x, y) \log \frac{1}{P_{XY}(x, y)} &= \sum_{x, y} P_{XY}(x, y) \log \frac{1}{P_{X|Y=y}(x)} + \sum_{x, y} P_{XY}(x, y) \log \frac{1}{P_Y(y)} \\ &= \sum_y P_Y(y) \sum_x P_{X|Y=y}(x|y) \log \frac{1}{P_{X|Y=y}(x)} + \sum_y P_Y(y) \log \frac{1}{P_Y(y)} \end{aligned}$$

which is equivalent to

Proposition 2 (The chain rule). $H(X, Y) = H(X|Y) + H(Y)$.

By treating multiple random variables as a single tuple of random variables, we can extend the chain rule to more variables. For example

$$\begin{aligned} H(X, Y|Z) &= H(X, Y, Z) - H(Z) = (H(X, Y, Z) - H(Y, Z)) + H(Y, Z) - H(Z) \\ &= H(X|Y, Z) + H(Y|Z). \end{aligned}$$

The *quantum conditional entropy* is *defined* to obey the chain rule.

Definition 3. For systems A and B :

1. The **quantum conditional entropy** of A given B is $H(A|B) := H(AB) - H(B)$.
2. The **coherent information** of A given B is $I(A|B) := -H(A|B)$.

Any identity between classical entropies derived using only the chain rule, can be derived in the same way for quantum entropies. For example

$$H(AB|C) = H(A|BC) + H(B|C).$$

A striking difference between classical and quantum conditional entropies is that while any classical conditional entropy is positive, quantum conditional entropy can be negative. For example, if $d_A = d_B = d$ and $\rho_{AB} = \phi_{AB}^+$ then $H(A|B)_\rho = H(AB)_\rho - H(B)_\rho = S(\phi_{AB}^+) - S(\mathbb{1}_B/d) = 0 - \log d$.

15.3 Mutual information and conditional mutual information

Given two random variables X and Y , the entropy gives us a way to quantify how much Y tells us about X : $H(X)$ is a way of quantifying our uncertainty about X . Prior to learning the value of Y , our expectation for the entropy we will assign to X if we *do* learn the value of Y , is $H(X|Y)$. Therefore, $I(X : Y) := H(X) - H(X|Y)$ is the expectation of the *reduction in entropy* of X which will occur if we learn the value of Y . This is a measure of how much learning Y would *inform* us about X . By the chain rule,

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

so $I(X : Y)$ is *symmetric* in the two random variables, and we call it the **mutual information** between X and Y .

Definition 4. For random variables X , Y and Z :

1. The mutual information between X and Y is

$$I(X : Y) := H(X) - H(X|Y). \quad (15.5)$$

2. The conditional mutual information between X and Y given $Z = z$ is

$$I(X : Y|Z = z) := H(X|Z = z) - H(X|Y, Z = z). \quad (15.6)$$

3. The conditional mutual information between X and Y given Z is

$$\begin{aligned} I(X : Y|Z) &:= \sum_z \Pr(Z = z) I(X : Y|Z = z) = H(X|Z) - H(X|Y, Z) \\ &= H(X|Z) + H(Y|Z) - H(X, Y|Z) \\ &= H(Y|Z) - H(Y|X, Z). \end{aligned}$$

Definition 5. For systems A , B and C :

1. The **quantum mutual information** between A and B is

$$I(A : B) := H(A) - H(A|B) = H(A) + H(B) - H(AB) = H(B) - H(B|A). \quad (15.7)$$

2. **quantum conditional mutual information** between A and B given C is

$$I(A : B|C) := H(A|C) - H(A|BC) = H(AC) + H(BC) - H(C) - H(ABC). \quad (15.8)$$

Ex. sheet 3 asks you to prove the **chain rule for conditional mutual information**:

Proposition 6. $I(X : Y, Z|W) = I(X : Z|W) + I(X : Y|Z, W)$.

Since this can be proved using the chain rule alone, it also holds for quantum entropies:

Proposition 7. $I(A : BC|D) = I(A : C|D) + I(A : B|CD)$.

15.4 Hypothesis testing and relative entropy

We now return briefly to a kind of simple state discrimination problem known as “quantum hypothesis testing”. Given a quantum system Q , let hypothesis H_0 be that the state of Q is ρ , and hypothesis H_1 be that the state of Q is σ . Suppose we measure a POVM $E : \{0, 1\} \rightarrow \mathcal{H}_Q$ obtaining a result \hat{X} which is supposed to identify which hypothesis is true. Since this is a binary POVM, if we set $E(0) = T$ then this determines $E(1) = \mathbb{1} - T$.

$$\text{The “type-I” error probability is } \Pr(\hat{X} = 1|H_0) = \alpha(T, \rho) := 1 - \text{Tr}T\rho, \text{ and} \quad (15.9)$$

$$\text{the “type-II” error probability is } \Pr(\hat{X} = 0|H_1) = \beta(T, \sigma) := \text{Tr}T\sigma. \quad (15.10)$$

In general, there is a trade-off between these two conditional probabilities. We use the following notation for the minimum value of β that can be attained subject to the requirement that $\alpha \leq \epsilon$.

Definition 8. $\beta_\epsilon(\rho||\sigma) := \min\{\beta(T, \sigma) : \alpha(T, \rho) \leq \epsilon, 0 \leq T \leq \mathbb{1}\}$.

Proposition 9 (Data processing inequality for β_ϵ). For all $\epsilon \in [0, 1]$, states ρ_A and σ_A , and operations $\mathcal{N}^{B \leftarrow A}$, $\beta_\epsilon(\mathcal{N}\rho_A||\mathcal{N}\sigma_A) \geq \beta_\epsilon(\rho_A||\sigma_A)$.

Proof. Because the linear map \mathcal{N} is completely positive and trace preserving, its adjoint \mathcal{N}^\dagger is completely positive and unital (identity preserving). So if $0 \leq T \leq \mathbb{1}$ then $0 \leq \mathcal{N}^\dagger T \leq \mathcal{N}^\dagger \mathbb{1} = \mathbb{1}$ and therefore

$$\beta_\epsilon(\mathcal{N}\rho||\mathcal{N}\sigma) = \min\{\beta(T, \mathcal{N}\sigma) : \alpha(T, \mathcal{N}\rho) \leq \epsilon, 0 \leq T \leq \mathbb{1}\} \quad (15.11)$$

$$= \min\{\beta(\mathcal{N}^\dagger T, \sigma) : \alpha(\mathcal{N}^\dagger T, \rho) \leq \epsilon, 0 \leq T \leq \mathbb{1}\} \geq \beta_\epsilon(\rho||\sigma). \quad (15.12)$$

□

Intuitively, if we are given n systems in the state ρ or n systems in the state σ then larger values of n should make it easier to distinguish between the two situations. The next theorem tells us that $\beta_\epsilon(\rho^{\otimes n}||\sigma^{\otimes n})$ exhibits exponential decay as n increases, and gives us the rate of the decay.

Theorem 10 (Quantum Stein’s lemma). For all states ρ, σ of a given system

$$\forall \epsilon \in (0, 1), \lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_\epsilon(\rho^{\otimes n}||\sigma^{\otimes n}) = -D(\rho||\sigma), \quad (15.13)$$

where $D(\rho||\sigma)$ is the *quantum relative entropy* between ρ and σ ...

Definition 11 (Quantum relative entropy).

$$D(\rho||\sigma) := \begin{cases} -S(\rho) - \text{Tr}\rho \log \sigma & \text{if } \text{supp}(\rho) \subseteq \text{supp}(\sigma), \\ +\infty & \text{otherwise.} \end{cases} \quad (15.14)$$

In $\text{Tr}\rho \log \sigma$ we regard both ρ and $\log \sigma$ as operators on $\text{supp}(\sigma)$.

Note that the relative entropy is not symmetric in its arguments and does not obey the triangle inequality.

Proposition 12. Density operators ρ and σ on \mathcal{H} commute if and only if there are eigendecompositions $\rho = \sum_{i=1}^{\dim(\mathcal{H})} p(i)|\alpha_i\rangle\langle\alpha_i|$ and $\sigma = \sum_{i=1}^{\dim(\mathcal{H})} q(i)|\alpha_i\rangle\langle\alpha_i|$. In this case, $D(\rho\|\sigma) = D(p\|q)$, where $D(p\|q)$ is the *classical relative entropy* between the distributions p and q ...

Definition 13 (Classical relative entropy).

$$D(p\|q) := \begin{cases} -S(p) - \sum_{x \in \text{supp}(q)} p(x) \log q(x) & \text{if } \text{supp}(p) \subseteq \text{supp}(q), \\ +\infty & \text{otherwise.} \end{cases} \quad (15.15)$$

Theorem 14 (Gibbs' inequality). For any two probability distributions p and q on a finite set \mathcal{A} , $D(p\|q) \geq 0$ with equality iff $p = q$.

Proof. That $D(p\|q) \geq 0$ already follows from the operational meaning given to $D(p\|q)$ by Stein's lemma. Here is a direct proof which gives us the equality condition. If $\text{supp}(p) \not\subseteq \text{supp}(q)$ then $D(p\|q) = \infty$ and the result holds, so let's now assume that $\text{supp}(p) \subseteq \text{supp}(q)$. $D(p\|q) = \sum_{x \in \text{supp}(p)} p(x) \log \frac{p(x)}{q(x)} = -\frac{1}{\ln(2)} \sum_{x \in \text{supp}(p)} p(x) \ln \frac{q(x)}{p(x)}$. The line $y = x - 1$ is tangent to the graph $y = \ln x$ at $x = 1$ so from the strict concavity of \ln it follows that $\ln x \leq x - 1$ with equality iff $x = 1$. Therefore,

$$-D(p\|q) \leq \frac{1}{\ln(2)} \sum_{x \in \text{supp}(p)} p(x) \left(\frac{q(x)}{p(x)} - 1 \right) \leq 0.$$

Since the $p(x)$ in the sum are strictly positive, the first inequality is an equality if and only if $q(x)/p(x) = 1$ for all $x \in \text{supp}(p)$, which is true iff $q(x) = p(x)$ for all x , in which case the second inequality is also satisfied. \square

Theorem 15 (Data processing inequality for D). For all states $\rho_{\mathcal{A}}$ and $\sigma_{\mathcal{A}}$, and operations $\mathcal{N}^{\mathcal{B} \leftarrow \mathcal{A}}$, $D(\mathcal{N}\rho_{\mathcal{A}}\|\mathcal{N}\sigma_{\mathcal{A}}) \leq D(\rho_{\mathcal{A}}\|\sigma_{\mathcal{A}})$.

Proof.

$$D(\mathcal{N}\rho\|\mathcal{N}\sigma) = - \lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_{\epsilon}((\mathcal{N}\rho)^{\otimes n} \| (\mathcal{N}\sigma)^{\otimes n}) \quad (15.16)$$

$$= - \lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_{\epsilon}(\mathcal{N}^{\otimes n} \rho^{\otimes n} \| \mathcal{N}^{\otimes n} \sigma^{\otimes n}) \quad (15.17)$$

$$\leq - \lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_{\epsilon}(\rho^{\otimes n} \| \sigma^{\otimes n}) = D(\rho\|\sigma). \quad (15.18)$$

The first and last equalities are by Quantum Stein's lemma. The inequality is by the data processing inequality for β_{ϵ} , the fact that \log is an increasing function, and an elementary fact about limits. \square

Theorem 16 (Klein's inequality). $D(\rho\|\sigma) \geq 0$ with equality iff $\rho = \sigma$.

Proof. Suppose that $X = 0$ with probability 1/2 and $X = 1$ with probability 1/2, and suppose that the state of \mathcal{Q} is ρ when $X = 0$ and the state is σ when $X = 1$. Suppose we measure a POVM E to obtain an estimate \hat{X} of X and then prepare \mathcal{Q} in the computational basis state $|\hat{X}\rangle\langle\hat{X}|$. This corresponds to applying the operation

$$\mathcal{N} : \tau \mapsto |0\rangle\langle 0| \text{Tr} E(0)\tau + |1\rangle\langle 1| \text{Tr} E(1)\tau. \quad (15.19)$$

We choose the POVM which maximises $\Pr(\hat{X} = X)$.

Letting $p(i) = \Pr(\hat{X} = i|X = 0)$ and $q(i) = \Pr(\hat{X} = i|X = 1)$, $\mathcal{N}\rho = p(0)|0\rangle\langle 0| + p(1)|1\rangle\langle 1|$ and $\mathcal{N}\sigma = q(0)|0\rangle\langle 0| + q(1)|1\rangle\langle 1|$. By the data processing inequality for quantum relative entropy, Proposition 12, and Gibbs' inequality

$$D(\rho\|\sigma) \geq D(\mathcal{N}\rho\|\mathcal{N}\sigma) = D(p\|q) \geq 0 \quad (15.20)$$

so we require $p = q$ for $D(\rho\|\sigma) = 0$ to hold. The Holevo-Helstrom theorem tells us that

$$\Pr(\hat{X} = X) = \frac{1}{2}(p(0) + q(1)) = \frac{1}{2}(1 + q(1) - p(1)) = \frac{1}{2}\left(1 + \frac{1}{2}\|\sigma - \rho\|_1\right) \quad (15.21)$$

so $q(1) - p(1) = \frac{1}{2}\|\sigma - \rho\|_1$, and $q = p$ iff $\|\sigma - \rho\|_1 = 0$ iff $\sigma = \rho$. \square

15.5 Entropic inequalities

Example sheet 3 asks you to show

Proposition 17. $I(X : Y) = D(P_{XY}\|P_X P_Y)$.

Therefore, Gibb's inequality tells us that $I(X : Y) \geq 0$ with equality iff $P_{XY} = P_X P_Y$. From this and the definitions of the classical quantities in section 15.3 we obtain

Proposition 18. For any random variables X, Y and Z (taking values in finite sets):

1. $I(X : Y) \geq 0$ with equality iff X and Y are independent, which means

$$\forall x, y : P_{XY}(x, y) = P_X(x)P_Y(y).$$

2. $I(X : Y|Z = z) \geq 0$ with equality iff X and Y are conditionally independent given $Z = z$, which means

$$\forall x, y : P_{XY|Z=z}(x, y) = P_{X|Z=z}(x)P_{Y|Z=z}(y).$$

3. $I(X : Y|Z) \geq 0$ with equality iff X and Y are conditionally independent given Z , which means

$$\forall z \in \text{supp}(P_Z), x, y : P_{XY|Z}(x, y|z) = P_{X|Z}(x|z)P_{Y|Z}(y|z).$$

Example sheet 3 also asks you to show

Proposition 19. $I(A : B)_\rho := D(\rho_{AB}\|\rho_A \otimes \rho_B)$.

So Klein's inequality gives us

Proposition 20. For any state ρ_{AB} , $I(A : B)_\rho \geq 0$ with equality iff $\rho_{AB} = \rho_A \otimes \rho_B$.

Quantum *conditional* mutual information is also positive, a fact known as the **strong subadditivity** of von Neumann entropy. Proving this from scratch is much more involved than the classical case, but it is a fairly easy consequence of the data processing inequality for quantum relative entropy.

Theorem 21. For any state ρ_{ABC} , $I(A : B|C) \geq 0$.

Proof. By the chain rule and Proposition 19,

$$\begin{aligned} I(A : B|C) &= I(A : BC) - I(A : C) \\ &= D(\rho_{ABC} \| \rho_A \otimes \rho_{BC}) - D(\rho_{AC} \| \rho_A \otimes \rho_C) \\ &= D(\rho_{ABC} \| \rho_A \otimes \rho_{BC}) - D(\text{Tr}_B \rho_{ABC} \| \text{Tr}_B \rho_A \otimes \rho_{BC}) \\ &\geq D(\rho_{ABC} \| \rho_A \otimes \rho_{BC}) - D(\rho_{ABC} \| \rho_A \otimes \rho_{BC}) = 0, \end{aligned}$$

where the inequality is the data processing inequality for quantum relative entropy. \square

Proposition 22. If ρ_{AB} is a pure state, then $H(A)_\rho = H(B)_\rho$.

Proof. From the existence of a Schmidt decomposition for ρ_{AB} we know that ρ_A and ρ_B have the same non-zero eigenvalues with the same multiplicities, and therefore have the same von Neumann entropy. \square

Theorem 23. Data processing inequality for the coherent information:

If $\sigma_{AC} = \mathcal{N}^{C \leftarrow B} \rho_{AB}$ for some operation $\mathcal{N}^{C \leftarrow B}$ then $I(A)B)_\rho \geq I(A)C)_\sigma$.

Proof. Let ρ_{RAB} be a purification of ρ_{AB} , and let $\mathcal{N}^{C \leftarrow B} X_B = \text{Tr}_E V X_B V^\dagger$ be a Stinespring representation for $\mathcal{N}^{C \leftarrow B}$ (so V is an isometry in $\mathcal{L}(\mathcal{H}_B, \mathcal{H}_C \otimes \mathcal{H}_E)$). Defining

$$\psi_{\text{RACE}} = \mathbb{1}_{RA} \otimes V \rho_{RAB} \mathbb{1}_{RA} \otimes V^\dagger,$$

we have $\text{Tr}_E \psi_{\text{RACE}} = \sigma_{AC}$, so let's write $\sigma_{\text{RACE}} := \psi_{\text{RACE}}$. Note that σ_{RACE} is a pure state. Now, the coherent information before the operation (the 'data processing') is

$$I(A)B)_\rho = H(B)_\rho - H(AB)_\rho \stackrel{(a)}{=} H(RA)_\rho - H(R)_\rho \stackrel{(b)}{=} H(RA)_\sigma - H(R)_\sigma = H(A|R)_\sigma, \quad (15.22)$$

where (a) is by two applications of Proposition 22 to the pure state ρ_{RAB} , (b) is because $\rho_{RA} = \sigma_{RA}$, and the other equalities are true by definition. After the operation,

$$\begin{aligned} I(A)C)_\sigma &= H(C)_\sigma - H(AC)_\sigma \stackrel{(c)}{=} H(RAE)_\sigma - H(RE)_\sigma \\ &= H(AE|R)_\sigma + H(R)_\sigma - (H(E|R)_\sigma + H(R)_\sigma) = H(AE|R)_\sigma - H(E|R)_\sigma, \end{aligned} \quad (15.23)$$

where (c) is by two applications of Proposition 22 to the pure state ρ_{RACE} and the other equalities are definitions. Subtracting equation (15.23) from (15.22) yields

$$I(A)B)_\rho - I(A)C)_\sigma = H(A|R)_\sigma + H(E|R)_\sigma - H(AE|R)_\sigma = I(A : E|R) \geq 0$$

by definition and positivity of QCMI. \square

Theorem 24. Data processing inequality for the quantum mutual information:

If $\sigma_{A'B'} = \mathcal{N}^{A' \leftarrow A} \otimes \mathcal{M}^{B' \leftarrow B} \rho_{AB}$ where $\mathcal{N}^{A' \leftarrow A}$ and $\mathcal{M}^{B' \leftarrow B}$ are operations then

$$I(A : B)_\rho \geq I(A' : B')_\sigma.$$

Proof. With $\rho'_{AB'} := \text{id}^{A \leftarrow A} \otimes \mathcal{M}^{B' \leftarrow B} \rho_{AB}$, $\sigma_{A'B'} = \mathcal{N}^{A' \leftarrow A} \otimes \text{id}^{B' \leftarrow B} \rho'_{AB'}$. Using $\rho'_A = \rho_A$ and $\rho'_{B'} = \sigma_{B'}$, the relationships between mutual information and coherent information, and applying the DPI for coherent information twice

$$\begin{aligned} I(A : B)_\rho &= H(A)_\rho + I(A)B)_\rho \geq H(A)_{\rho'} + I(A)B')_{\rho'} \\ &= I(A : B')_{\rho'} = H(B')_{\rho'} + I(B')A)_{\rho'} \geq H(B')_\sigma + I(B')A)_\sigma = I(A' : B')_\sigma. \end{aligned}$$

\square

15.6 Fano's inequality

If $\mathcal{A}_X = \{0, 1\}$ then $H(X) = h(P_X(1))$ where h is the **binary entropy function**:

Definition 25. $h(\lambda) := (1 - \lambda) \log \frac{1}{1-\lambda} + \lambda \log \frac{1}{\lambda}$.

Proposition 26 (Fano's inequality). Given two random variables M and \hat{M} which both take values in a finite set \mathcal{A} , we have

$$H(M|\hat{M}) \leq \Pr(\hat{M} \neq M) \log(|\mathcal{A}| - 1) + h(\Pr(\hat{M} \neq M)).$$

Proof. Let E be a binary RV which is 1 if $\hat{M} \neq M$ and 0 if $\hat{M} = M$. Using the chain rule and $H(E|M, \hat{M}) = 0$ (because E is a function of M and \hat{M}):

$$\begin{aligned} H(M|\hat{M}) &= H(M|\hat{M}, E) + H(E|\hat{M}) - H(E|M, \hat{M}) \\ &= H(M|\hat{M}, E = 1) \Pr(E = 1) + H(M|\hat{M}, E = 0) \Pr(E = 0) + H(E|\hat{M}) \end{aligned}$$

$H(M|\hat{M}, E = 0) = 0$ because, conditioned on $E = 0$ (i.e. $\hat{M} = M$) M is a function of \hat{M} . $H(M|\hat{M}, E = 1) \leq \log(|\mathcal{A}| - 1)$ because, for all $a \in \mathcal{A}$, the support of $P_{M|\hat{M}=a, E=1}$ is no larger than $|\mathcal{A}| - 1$. Finally, by positivity of mutual information, $H(E|\hat{M}) \leq H(E) = h(\Pr(\hat{M} \neq M))$. \square