

## Review

- $\frac{d}{dx}(x \ln x) = \ln x + 1 = \ln ex$
- $\frac{d}{dx}(x \log x) = \log x + \frac{1}{\ln 2} = \log ex$
- $\frac{d}{dx}(-x \log x) = -\log x - \frac{1}{\ln 2}$
- $h(p) = -p \log p - (1-p) \log(1-p)$ 
  - $\frac{dh}{dp}(p) = \log \frac{1-p}{p}$
  - $\frac{d}{dx}h(f(x)) = \left(\frac{df}{dx}(x)\right) \log \frac{1-f(x)}{f(x)}$
  - $2^{h(p)} = p^{-p} (1-p)^{-(1-p)}$
  - $2^{-h(p)} = p^p (1-p)^{1-p}$
  - $h\left(\frac{1}{b}\right) = \log b - \frac{b-1}{b} \log(b-1)$
- **Probability transition matrix:**  $[Q(y|x)]_{r,c} = \Pr[y=c|x=r]$ . (The entry in the  $x^{\text{th}}$  row and the  $y^{\text{th}}$  column denotes  $Q(y|x)$ )
  - Let input distribution be a row vector  $\vec{p}^T$ . Then the output distribution would be a row vector  $\vec{q}^T = \vec{p}^T Q$ .

**Markov String/Process:**  $X_1^n: p(x_1^n) = p(x_1) p(x_2|x_1) \cdots p(x_k|x_{k-1}) p(x_{k+1}|x_k) \cdots p(x_n|x_{n-1})$

- Ordered substring of Markov string is Markov  
For  $0 \leq n_1 < n_2 < \cdots < n_k \leq n$ ,  $(X_{n_1}, X_{n_2}, \dots, X_{n_k})$  is also a Markov string.
- Given  $X_k$  (the present), we have  $X_1^{k-1}$  (the past) and  $X_{k+1}^n$  (the future) are independent:
  - $H(X_1^{k-1}, X_{k+1}^n | X_k) = H(X_1^{k-1} | X_k) + H(X_{k+1}^n | X_k)$
  - $I(X_1^{k-1}; X_{k+1}^n | X_k) = 0$
- Given  $X_k$ , we have  $X_{k-1}$  and  $X_{k+1}$  are independent:
  - $p(x_{k-1}, x_{k+1} | x_k) = p(x_{k-1} | x_k) p(x_{k+1} | x_k)$
- For  $1 \leq k < n$ , and  $k+1 < m \leq n$ ,
  - $H(X_{k+1}^m | X_k) = \sum_{i=k+1}^m H(X_i | X_k)$

- $p(x_{k+1}^m | x_j^k) = p(x_{k+1}^m | x_k)$       •  $H(X_{k+1}^n | X_1^k) = H(X_{k+1}^n | X_k)$
- $H(X_0 | X_k)$  is increasing in  $k$ :  $H(X_0 | X_{n+1}) \geq H(X_0 | X_n)$
- $I(X_0; X_k)$  is decreasing in  $k$ :  $I(X_0; X_n) \geq I(X_0; X_{n+1})$
- Markov 3-string:  $X \text{---} \ominus \text{---} Y \text{---} \ominus \text{---} Z$ :  $p(x, y, z) = p(x)p(y|x)p(z|y)$ 
  - Also Markov in the reverse direction:  $p(z, y, x) = p(x|y, z)p(y, z) = p(x|y)p(y|z)p(z)$ .
  - $I(Z; X | Y) = I(X; Z | Y) = 0$ .
  - **Data processing theorem**:  $I(X; Y) \geq I(X; Z)$ ,  $I(Z; Y) \geq I(Z; X)$ ; hence, closer  $\Rightarrow$  more  $I$ .
  - $I(X; Y) \geq I(X; Y | Z)$
- $\underline{U} \text{---} \ominus \text{---} \underline{X} \text{---} \ominus \text{---} \underline{Y} \text{---} \ominus \text{---} \underline{V} \Rightarrow I(\underline{X}; \underline{Y}) \geq I(\underline{U}, \underline{V})$ .
- Stationary Markov process:  $H(X_n)$  is constant. (by stationarity).  $H(X_n | X_1)$  increases with  $n$ .

## Convexity:

- $H(Y)$  is a concave  $\cap$  function of  $p(x)$  for fixed  $Q(y|x)$ .
- $H(Y|X)$  is a linear function of  $p(x)$  for fixed  $Q(y|x)$ .
- $I(X; Y)$  is a continuous concave  $\cap$  function of  $p(x)$  for fixed  $Q(y|x)$ .
- $I(X; Y)$  is a convex  $\cup$  function of  $Q(y|x)$  for fixed  $p(x)$ .

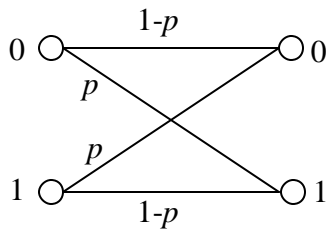
## $I(X; Y)$

- $I(X; Y) = 0$  if  $|\mathcal{X}| = 1$  or  $|\mathcal{Y}| = 1$ .
- To find  $I(p, Q)$ , first find  $q(y)$ . Then, find  $H(Y)$ . Next, find  $H(Y|X) = \sum_x p(x)H(Y|X=x)$ . Finally,  $I(X; Y) = I(p, Q) = H(Y) - H(Y|X)$ .
- **Parallel d.m.c. channel.** Let  $Y_1^n$  be the result of passing  $X_1^n$  through a d.m.c. ( $n$  use.)
  - $Q(y_1^n | x_1^n) = \prod_{i=1}^n Q_i(y_i | x_i)$ .
  - $H(Y_1^n | X_1^n) = \sum_{i=1}^n H(Y_i | X_i)$
  - $\forall p(x_1^n) \quad I(X_1^n; Y_1^n) \leq \sum_{i=1}^n I(X_i; Y_i)$  with equality if  $p(x_1^n) = \prod_{i=1}^n p_i(x_i)$ .
- Independent source:  $U_1^L$  has independent components  $\Rightarrow I(U_1^L; V_1^L) \geq \sum_{\ell=1}^L I(U_\ell; V_\ell)$ .

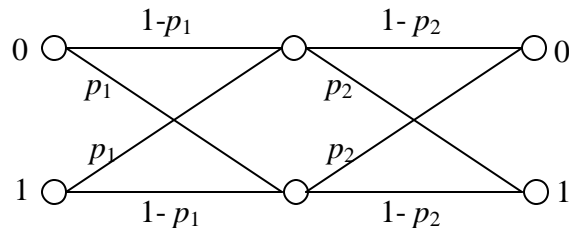
- Let  $U, V$  discrete random variables.  $M = |\mathcal{U}| = |\mathcal{V}|$ .  $\mathcal{U} = \mathcal{V} = \{0, 1, \dots, M-1\}$ .
  - Fano Inequality:**  $H(U|V) \leq h(P_e) + P_e \log(M-1)$ .
    - Note:  $P_e = 0 \Rightarrow H(U|V) = 0$ .
  - Extended Fano inequality:** Let  $U_1^L, V_1^L \in \mathcal{U}^L = \mathcal{V}^L$ .  $\frac{H(U_1^L|V_1^L)}{L} \leq h(\bar{P}_e) + \bar{P}_e \log(M-1)$

## BSC( $p$ )

- Binary symmetric channel (BSC) with crossover probability  $p$ .  $x \longrightarrow \boxed{\text{BSC}(p)} \longrightarrow \underline{y}$ .



- $\begin{bmatrix} 1-p & p \\ p & 1-p \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1-2p & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}^{-1}$ .
- $p$  also = probability of error.
- $H(Y|X) = h(p)$  regardless of  $\{p(x)\}$ . (by the “symmetric”)
- $C = 1-h(p)$ .
- Two BSC’s in series is a BSC with transition probability  $p^{(2)} = p_1q_2 + p_2q_1$ .  $q^{(2)} = p_1p_2 + q_1q_2$ .



$$Q^{(2)} = Q_1 Q_2 = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} (1-2p_1)(1-2p_2) & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}^{-1}$$

- $n$  BSC’s in series is a BSC with transition probability  $p^{(n)} = \frac{1 - \prod_{i=1}^n (1-2p_i)}{2}$ .

$$Q^{(n)} = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \prod_{i=1}^n (1-2p_i) & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}^{-1}$$

- Two identical BSC's in series is a BSC with transition probability  $p^{(2)} = 2pq$ .
- $n$  identical BSC's in series is a BSC with transition probability  $p^{(n)} = \frac{1 - (1 - 2p)^n}{2}$ .

## Binary $n$ -Sphere

- Binary  $n$ -cube: geometric representation of  $\{0,1\}^n$ .  $\underline{x} = (x_1, \dots, x_n) \in \{0,1\}^n$ .
- Combinatoric fact:
  - $\binom{n}{j}$  is increasing for  $j \leq \left\lfloor \frac{n}{2} \right\rfloor$  and max when  $j = \left\lfloor \frac{n}{2} \right\rfloor$ .
  - **Stirling's approximation:**  $\sqrt{2\pi} n^n n^{\frac{1}{2}} e^{-n+\frac{1}{12n+1}} < n! < \sqrt{2\pi} n^n n^{\frac{1}{2}} e^{-n+\frac{1}{12n}} \Rightarrow \ln n! = n \ln n - n + o(n)$  as  $n \rightarrow \infty$ .  $\left( \lim_{n \rightarrow \infty} \frac{o(n)}{n} = 0 \right)$ .
  - $\log_D n! = n \log n - \frac{1}{\ln D} n + o(n) = n \log n + an + o(n)$
- The **hamming sphere** of radius  $r$  around  $\underline{x}$  is  $S_r(\underline{x}) = \{ \underline{y} \in \{0,1\}^n ; d_H(\underline{x}, \underline{y}) \leq r \}$ .
  - $|S_r(\underline{x})| = \text{Volume of radius } r \text{ sphere} = \sum_{i=0}^r \binom{n}{i}$ .
  - Let  $S_k = S_k(\bar{0}) \subset \{0,1\}^n$ .
  - $\binom{n}{k} \leq |S_k| \leq (k+1) \binom{n}{k}$
  - Let  $k = \mathbf{a}n$  where  $0 < \mathbf{a} \leq \frac{1}{2}$ . Then,  $\lim_{n \rightarrow \infty} \frac{1}{n} \log_{\text{volume}} |S_{\mathbf{a}n}| = \lim_{n \rightarrow \infty} \frac{1}{n} \log \sum_{i=0}^{\mathbf{a}n} \binom{n}{i} = \lim_{n \rightarrow \infty} \frac{1}{n} \log \binom{n}{\mathbf{a}n} = H(\mathbf{a})$ : The rate of growth of exponential of the volume is the same as that of the surface.  $|S_{\mathbf{a}n}| \sim 2^{nH(\mathbf{a})}$ .
  - Multinomial extension: for  $\sum_{i=1}^M \mathbf{a}_i n = 1$ ,  $\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{n!}{\prod_{i=1}^M (\mathbf{a}_i n)!} = H(\mathbf{a}_1, \dots, \mathbf{a}_n)$ .
- $d_H(\underline{x}, \underline{y}) = \|\underline{x} - \underline{y}\| = |\{k : 1 \leq k \leq n, x_k \neq y_k\}|$
- Given  $\underline{x} \in \{0,1\}^n$ ,  $\binom{n}{i}$  = the number of sequences  $\underline{y} \in \{0,1\}^n$  with  $d_H(\underline{x}, \underline{y}) = i$ .

## Coding

- **Hamming  $(n,k)$  code** consists of  $2^k$  elements in  $\{0,1\}^n$
- **$e$ -error correcting perfect code** : every  $\underline{x} \in \{0,1\}^n$  is in  $e$ -sphere around some codeword. These  $e$ -spheres don't overlap.
  - The only binary perfect codes are
    - Hamming
    - Golay code ( $e = 3$ )
    - Repetition code ( $e = m$ ) where  $n = 2m+1$  (w/ majority vote)
- **$e$ -error correcting perfect hamming code** :  $\left( \sum_{i=0}^e \binom{n}{i} \right) 2^{n-k} = 2^n$ , or equivalently,  $\sum_{i=0}^e \binom{n}{i} = 2^k$ .
- **$e$ -error correcting sphere-packed code** : every  $\underline{x} \in \{0,1\}^n$  is in  $(e+1)$ -sphere around some codeword. These  $e$ -spheres don't overlap.
- Hamming  $(7,4)$  code is
  - not unique unless require  $\bar{0}$  in the code.
  - $S_1(\underline{x}) \cap S_1(\underline{y}) = \emptyset$  if  $\underline{x}, \underline{y} \in \mathcal{C}$  and  $\underline{x} \neq \underline{y}$
  - 1-error correcting perfect code:  $\left( \binom{7}{0} + \binom{7}{1} \right) 2^{7-4} = 2^7$ .
- Hamming  $(23,11)$  code
  - 3-error correcting perfect code:  $\left( \sum_{i=0}^3 \binom{23}{i} \right) 2^{23-11} = 2^{23}$ .
- **Rate of an  $(n,k)$  code** is  $R = \frac{k}{n}$ .  $\Rightarrow$  rate  $\frac{k}{n}$  code.
- $\forall R$   $0 < R < 1$ ,  $\exists n(R)$  such that no sphere-packed codes of blocklength  $n > n(R)$  exist.
- Geometric arguments
  - $\forall \epsilon > 0$  as  $n \rightarrow \infty$ , received word when  $\underline{x}$  is sent falls in the spherical shell of width  $\pm n\epsilon$  or normalized width  $\pm \epsilon$ .
  - $\exists R^* > 0 = \max$  of  $R > 0$  such that  $2^{nR}$  words can be put into  $\{0,1\}^n$  with sphere of radius  $np$  around them not overlapping.
  - **Sphere hardening**: For any  $\epsilon > 0$ , as  $n \rightarrow \infty$ ,  $\underline{y} \in S = S_{np+n\epsilon}(\underline{x}) \setminus S_{np-n\epsilon}(\underline{x})$  = the sphere of radius  $np$  with width  $2n\epsilon$  with probability 1.
  - For BSC( $p$ ),  $R^* > C = 1 - H(p) = 1 + p \log p + (1-p) \log(1-p)$  ??
- [Block Coding](#)

- Block code:  $\mathcal{C} \subset \{0,1\}^n$ . Blocklength  $n$ .
- Size:  $|\mathcal{C}| = \#$  of codewords.
- **Rate**  $R_c = \boxed{R = \frac{\log|\mathcal{C}|}{n}}$  [(info) bits per channel use] if messages are a priori equiprobable.

### Capacity for Discrete memoryless channel.

- An  $(n, R, I)$  code for a discrete memoryless channel with input alphabet  $\mathcal{X}$ , and output alphabet  $\mathcal{Y}$ , is a collection  $\mathcal{C}$  of  $2^{nR}$  codewords each belonging to  $\mathcal{X}^n$ . ( $2^{nR} \leq |\mathcal{X}|^n$ ).  $I$  is the probability of error.
- (block length, # codewords, minPr[error]) code.
- Some idea: For BSC( $p$ ), we have sphere of radius  $np$  with volume (or surface)  $\sim 2^{nH(p)}$ . Whole space has  $2^n$  sequences. So, if can build perfect code for which  $np$ -spheres around words don't overlap, then such code would have  $\frac{2^n}{2^{nH(p)}} = 2^{n(1-H(p))}$  codewords. Code's rate would be  $\frac{1}{n} \log_2 2^{n(1-H(p))} = 1 - H(p)$ .
- **Operational definition of capacity**: capacity of a discrete memoryless channel is the sup of all rates  $R$  such that there is a sequence of  $(n, 2^{nR}, I)$  codes for which  $I \rightarrow 0$  as  $n \rightarrow \infty$ .
- Consider discrete memoryless channel with transition matrix  $[Q(y|x)], x \in \mathcal{X}$  and  $y \in \mathcal{Y}, |\mathcal{X}|, |\mathcal{Y}| < \infty$ . Let  $I(p, Q) = I(X; Y)$  be the average mutual information between channel input and channel output when input random variable has p.m.f.  $\{p(x), x \in \mathcal{X}\}$ . Then,  $\boxed{C = \max_p I(p, Q)}$  [bits per channel use].
- $C$  is unique, but  $\operatorname{argmax}_p I(p, Q)$  may not be unique.

### Discrete Memoryless Channel (DMC)

- $\underline{x} = (x_1, \dots, x_n), \underline{y} = (y_1, \dots, y_n)$ .  $\underline{x} \longrightarrow \boxed{\text{Channel } Q_{y|x}(\underline{y}|\underline{x})} \longrightarrow \underline{y}$ .
- $Q_{y|x}(\underline{y}|\underline{x}) = \prod_{k=1}^n Q(y_k|x_k)$ .  $\{Q(y|x), x \in X, y \in Y\}$  is fixed.
- **Probability transition matrix**:  $[Q(y|x)]_{r,c} = \Pr[y=c|x=r]$ . (The entry in the  $x^{\text{th}}$  row and the  $y^{\text{th}}$  column denotes  $Q(y|x)$ .)

$$X \begin{bmatrix} \ddots & \vdots & \ddots \\ \cdots & Q(Y=y|X=x) & \cdots \\ \ddots & \vdots & \ddots \end{bmatrix} Y$$

- $P(x_1^n, y_1^n) = p(x_1^n) \prod_{i=1}^n Q(y_i | x_i)$
- Let  $N = \{1, \dots, n\}$ , and  $A \subset M \subset N$ , then  $p(\bar{x}_M, \bar{y}_A) = p(x_M) \prod_{\ell \in A} Q(y_\ell | x_\ell)$ .
- $Q\left(y_i \left| \begin{matrix} x_{I_1}, y_{I_2} \\ i \in I_1 \quad i \notin I_2 \end{matrix} \right. \right) = Q(y_i | x_i)$ .

- $p(x_I, y_I) = p(x_I) \prod_{\ell \in I} p(y_\ell | x_\ell)$
- For  $i \in I \subset \{1, \dots, n\}$ ,  $p(x_I, y_i) = p(x_I) p(y_i | x_i)$
- $p\left(y_i \left| \begin{matrix} x_I \\ i \in I \end{matrix} \right. \right) = p(y_i | x_i)$
- $p(y_i | x_i, y_k) = p(y_i | x_i) \quad k \neq i$
- $p(y_1^k, x_1^n) = p(x_1^n) \prod_{i=1}^k p(y_i | x_i)$
- $p(y_i | y_1^{i-1}, x_1^n) = \frac{p(y_1^i, x_1^n)}{p(y_1^{i-1}, x_1^n)} = \frac{p(x_1^n) \prod_{k=1}^i p(y_k | x_k)}{p(x_1^n) \prod_{k=1}^{i-1} p(y_k | x_k)} = p(y_i | x_i)$ .

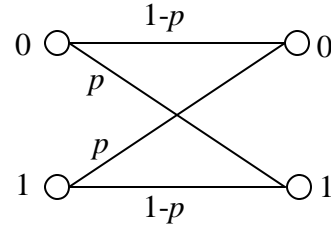
- $I(X_1^n; Y_1^n) \leq \sum_{i=1}^n I(X_i; Y_i)$

### D.M.C. with i.i.d. inputs

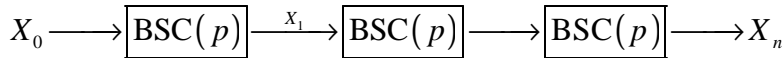
- Setup: Channel is d.m.c.  $\{Q(y|x)\}$ :  $Q(\bar{y}|\bar{x}) = \prod_{k=1}^n Q(y_k|x_k)$ . The input  $X_i$ 's to the channel is i.i.d.  
 $\Pr[X_1^n = x_1^n] = \prod_{k=1}^n p(x_k)$ .
- $(X_i, Y_i)$  is i.i.d. i.e.  $P(\bar{x}, \bar{y}) = \prod_{k=1}^n P(x_k, y_k) = \prod_{k=1}^n p(x_k) Q(y_k|x_k)$
- $Y_i$  is i.i.d. i.e.  $q(\bar{y}) = \prod_{k=1}^n q(y_k)$

## Capacity

- $0 \leq C \leq \min \{ \log |\mathcal{X}|, \log |\mathcal{Y}| \}$ .
- Any channel with only one input letter or only one output letter has zero capacity. ( $I(X;Y) \equiv 0$ ).
- BSC( $p$ ):  $C = 1 - H(p)$  is achieved when the input is uniform.
  - $H(Y|X) = H(p)$ .
  - $I(X;Y) = H(Y) - H(p)$ .



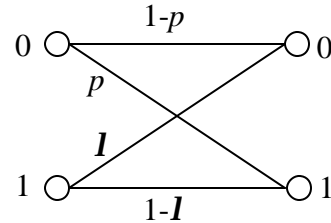
- Cascade of  $n$  identical BSC( $p$ ): is a BSC with transition probability  $p^{(n)} = \frac{1 - \prod_{i=1}^n (1 - 2p_i)}{2}$ .



- For  $0 < p < 1$ , because  $|1 - 2p| < 1$ ,  $\lim_{n \rightarrow \infty} p^{(n)} = \frac{1}{2}$ .  $\lim_{n \rightarrow \infty} C^{(n)} = 1 - H\left(\frac{1}{2}\right) = 0$ .  $0 \leq I(X_0; X_n) \leq C^{(n)} \Rightarrow \lim_{n \rightarrow \infty} I(X_0; X_n) = 0$  for any initial distribution.

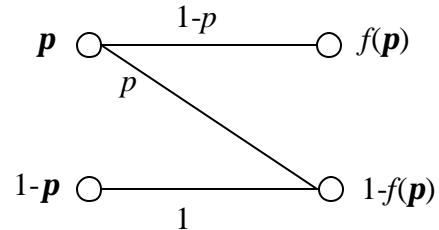
- Binary asymmetric channel:

- $f(\mathbf{p}_1) = (1 - p) + (p + \mathbf{1} - 1)\mathbf{p}_1$ .
- $H(Y) = h(f(\mathbf{p}_1))$ .
- $H(Y|X) = (1 - \mathbf{p}_1)h(p) + \mathbf{p}_1h(\mathbf{1})$ .
- For max  $I(X;Y)$ , need  $\log \frac{1 - f(\mathbf{p}_1)}{f(\mathbf{p}_1)} = \frac{h(\mathbf{1}) - h(p)}{p + \mathbf{1} - 1}$ .



- Z-channel:

- Let  $\mathbf{p}$  be the probability of the X that that channel introduces noise.
- $f(\mathbf{p}) = (1 - p)\mathbf{p}$ ;  $\log \frac{1 - f(\mathbf{p})}{f(\mathbf{p})} = \frac{h(p)}{1 - p}$ .
- $H(Y|X) = \mathbf{p}h(p)$ ;  $H(Y) = h(f(\mathbf{p}))$
- $C$  is achieved when





$$p = p^* = \frac{1}{(1-p) \left( 1 + 2^{\frac{h(p)}{1-p}} \right)} = \frac{1}{1-p + p^{\frac{p}{(1-p)}}}$$

$$C = h((1-p)p^*) - p^* h(p).$$

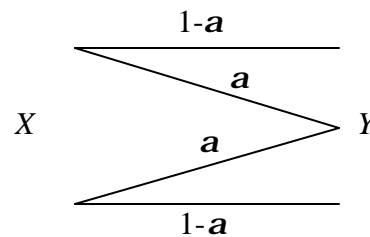
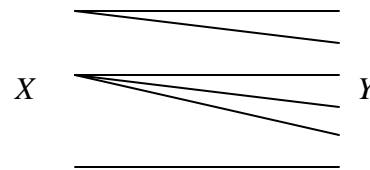
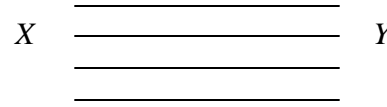
- Noiseless channel:  $C = \log|\mathcal{X}|$  is achieved by uniform input distribution.

- $H(Y|X) = 0$ .

- Noisy channel with nonoverlapping outputs:  $C = \log|\mathcal{X}|$  is achieved by uniform input distribution..

- $H(Y|X) = 0$ .

- Binary Erasure Channel:  $C = 1-a$  is achieved by uniform input distribution.



- Weakly symmetric channel:** 1) every row of the transition matrix are permutations of each other, i.e.,  $\{Q(y|i)\}$  are permutation of  $\{Q(y|j)\}$ , and 2) all the column sums  $\sum_x Q(y|x)$  are equal.

- $\sum_x Q(y|x) = \frac{|\mathcal{X}|}{|\mathcal{Y}|}$ .

- Uniform distribution on the input alphabet implies uniform distribution on output.
- $C = \log|\mathcal{Y}| - H(\text{rows of transition matrix})$  is achieved by a uniform distribution on the input alphabet.
- $\forall x H(Y|X = x) = H(\bar{r}^T)$ . Hence,  $H(Y|X) = H(\bar{r}^T)$  where  $\bar{r}^T$  is any row of the transition matrix..
- $I(X;Y) = H(Y) - H(\bar{r}^T)$

- Symmetric channels:** All the rows of the probability transition matrix  $[Q(y|x)]_{r,c} = \Pr[y=c|x=r]$  are permutations of each other and so are the columns, i.e. 1)  $\{Q(y|i)\}$  are just permutation of  $\{Q(y|j)\}$ , and 2)  $\{Q(i|x)\}$  are just permutation of  $\{Q(j|x)\}$ .

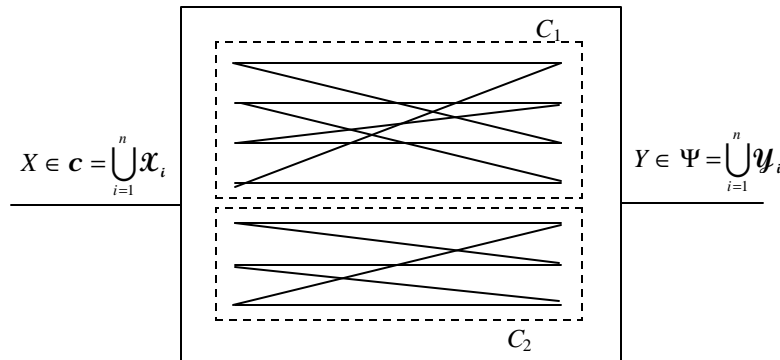
- Ex.  $\mathcal{X} = \mathcal{Z} = \{0,1,\dots,M-1\}$ .  $Y = (X + Z) \bmod M$ .

- Ex. BSC.

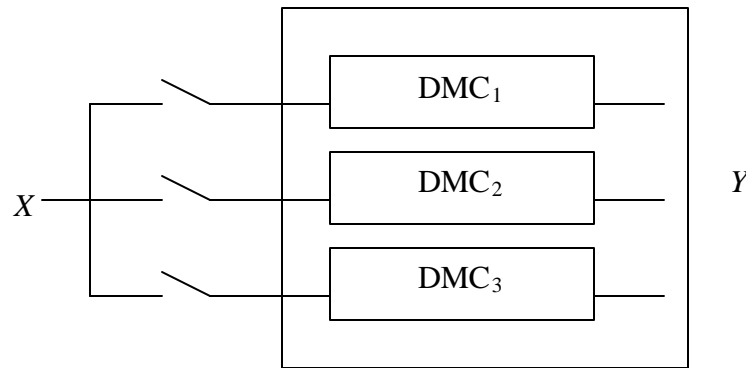
$\Rightarrow$  weakly symmetric. Hence,  $C = \log|\mathcal{Y}| - H(\text{rows of transition matrix})$ .

- **Sum channel**

- Consider  $N$  DMC's with disjoint input alphabets and disjoint output alphabets and capacities  $C_1, C_2, \dots, C_N$ . Call these DMC's as the **subchannels**. The associated **sum channel** has input and output alphabets that are the unions of those of the sub channel, and for each input  $x$ , the transition probabilities  $Q(\cdot|x)$  are the same as in the sub channel that has  $x$  in its input alphabet. In other words, the sum channel has all  $N$  subchannel available but only one subchannel may be used at any given time.



(Or can have common input alphabet but disjoint output alphabet, and selectable channel as shown below)



- Let  $\mathcal{X}_i$  and  $\mathcal{Y}_i$  be the input alphabet and the output alphabet of the  $i^{\text{th}}$  subchannel.  $\mathcal{C} = \bigcup_{i=1}^n \mathcal{X}_i$ ,  $\Psi = \bigcup_{i=1}^n \mathcal{Y}_i$ .  $\mathcal{X}_i \cap \mathcal{X}_j = \emptyset$ , and  $\mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset$  for  $i \neq j$ . Let  $w(n)$  be the probability that  $X \in \mathcal{X}_n$ , and  $p_n(x) = \Pr[X = x | X \in \mathcal{X}_n]$ . Then, for  $x \in \mathcal{X}_i$ ,  $p(x) = w(i) p_i(x)$ .
  - For  $x \in \mathcal{X}_i$  and  $y \in \mathcal{Y}_j$ ,  $Q(y|x) = Q_j(y|x) \mathbf{d}(j,i)$ .
- Consider each subchannel,
  - Let  $I_j(X;Y) = \sum_{x \in \mathcal{X}_j} p_j(x) \sum_{y \in \mathcal{Y}_j} Q_j(y|x) \log \frac{Q_j(y|x)}{\sum_{x \in \mathcal{X}_j} p_j(x) Q_j(y|x)}$ . Then,  $C_j = \max_{\{p_j(x), x \in \mathcal{X}_j\}} I_j(X;Y)$ .

- For  $y \in \mathcal{Y}_j$ ,  $q(y) = \sum_{x \in \mathcal{X}_j} w(j) p_j(x) Q_j(y|x)$ .

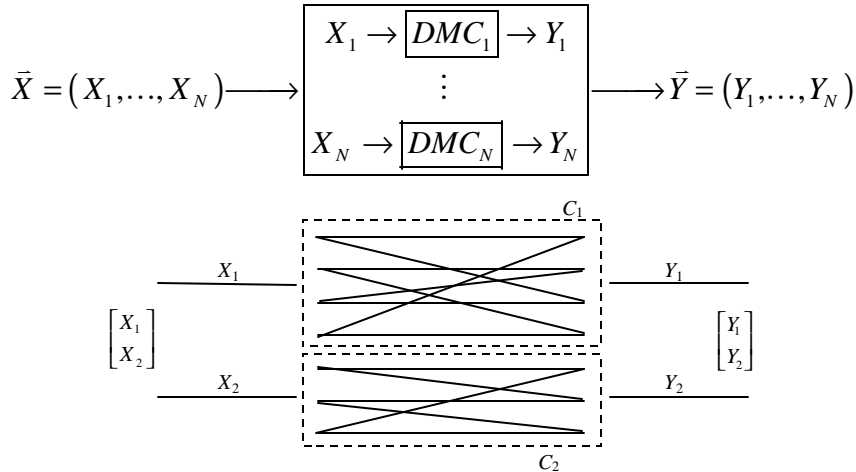
- $I(X;Y) = \sum_{x \in \mathcal{C}} p(x) \sum_{y \in \mathcal{Y}} Q(y|x) \log \frac{Q(Y|X)}{q(Y)} = \sum_{j=1}^N w(j) I_j(X;Y) - \sum_{j=1}^N w(j) \log w(j)$ .

- $C = \log \sum_{\ell=1}^N 2^{C_\ell}$  is achieved by  $w(j) = \frac{2^{C_j}}{\sum_{\ell=1}^N 2^{C_\ell}}$  and  $\{p_j(x), x \in \mathcal{X}_j\} = \arg \max_{\{p_j(x), x \in \mathcal{X}_j\}} I_j(X;Y)$ .

- Parallel channel:**

- Consider  $N$  DMC **subchannels** with capacities  $C_1, C_2, \dots, C_N$ . The subchannels are connected in parallel in the sense that once each unit of time an arbitrary symbol is transmitted and received over each subchannel. The output from each subchannel depends only on the input to that channel.

$$\left( Q(y_1^N | x_1^N) = \prod_{j=1}^N Q_j(y_j | x_j) \right).$$



- $C = \sum_{i=1}^n C_i$  is achieved by independent input  $p(x^n) = \prod_{i=1}^n \tilde{p}_i(x_i)$  where  $\{\tilde{p}_i(x)\}$  is the distribution that achieves  $C_i$  for the  $i^{\text{th}}$  subchannel.

- Cascade BSC's: A cascade of  $n$  identical BSCs each with transition probability  $p$  is equivalent to a single BSC with

- Iterative Calculation of  $C$ : The **Arimoto-Blahut Algorithm**.

- Given a DMC with transition probabilities  $Q(y|x)$  and any distribution  $p_0(x)$ .

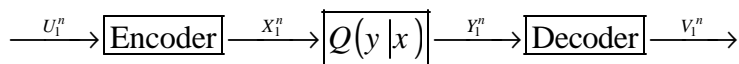
- Define a sequence  $p_r(x)$ ,  $r = 0, 1, \dots$  according to the iterative prescription:

$$q_r(y) = \sum_x p_r(x) Q(y|x) \cdot \log c_r(x) = \sum_y Q(y|x) \log \frac{Q(y|x)}{q_r(y)} \cdot p_{r+1}(x) = \frac{p_r(x) c_r(x)}{\sum_x p_r(x) c_r(x)}.$$

- 1) Set  $p_0 = [p_0(0), p_0(1), \dots]$  (row vector).  $Q = \begin{bmatrix} \Pr[0 \rightarrow 0] & \Pr[0 \rightarrow 1] \\ \Pr[1 \rightarrow 0] & \Pr[1 \rightarrow 1] \\ \vdots & \vdots \end{bmatrix}$ . Let  $Q(x)$  be the row  $x$  of  $Q$ .
- 2)  $q_r = p_r Q$ .
- 3) For each  $x$ ,  $\log c_r(x) = -H(\text{row } x \text{ of } Q) - \sum((\text{row } x \text{ of } Q) \cdot \log(q_r))$ .
- 4) For each  $x$ ,  $c_r(x) = 2^{\log c_r(x)}$ . Form  $c_r = c_r = [c_r(0), c_r(1), \dots]$ .
- 5)  $temp = \sum(p_r \cdot c_r)$ .
- 6)  $p_{r+1} = \frac{1}{temp}(p_r \cdot c_r)$ .
- 7) Repeat 2) – 6)
- $f(p_{r+1}, Q, \hat{P}_r) = \max_p f(p, Q, \hat{P}_r) \geq f(p_r, Q, \hat{P}_r)$
- $I(p_{r+1}(x), Q) \geq I(p_r(x), Q)$ ; thus,  $I(p_r, Q)$  is monotonic increasing with  $r$ .
- $\log\left(\sum_x p_r(x) c_r(x)\right) \leq C \leq \log\left(\max_x c_r(x)\right)$
- $\lim_{r \rightarrow \infty} I(p_r, Q) = C$  if  $\forall x p_0(x) > 0$

## System

- Feedback: all the received symbols  $Y_i$  are sent back immediately and noiselessly to the transmitter, which can then use them to decide which symbol to send next.
  - Feedback can help enormously in simplifying encoding and decoding. However, it can not increase the capacity of d.m.c.
- Source-Channel coding theorem: we can transmit a stationary ergodic source over a channel if and only if its entropy rate is less than the capacity of the channel.



- $P_e^{(n)} = \Pr[U_1^n \neq V_1^n] = \sum_{y_1^n} \sum_{u_1^n} p(u_1^n) Q(y_1^n | \text{Enc}(u_1^n)) I_{\{\text{Dec}(y_1^n) \neq u_1^n\}}$
- If  $U_1, U_2, \dots, U_n$  is a finite alphabet stochastic process that satisfies the AEP (ex. stationary ergodic source), then there exists a source channel code with  $P_e^{(n)} \rightarrow 0$  if source entropy rate  $H(\mathbf{u}) < C$ .
- For any stationary stochastic process, if  $H(\mathbf{u}) > C$ , the probability of error is bounded away from zero, and it is not possible to send the process over the channel with arbitrary low probability of error.

## Info transmission theorem with stationary source and DMC

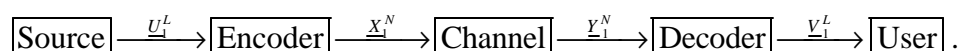
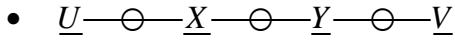


Fig 1.



• Let  $\mathcal{U} = \mathcal{V} = \{0, 1, \dots, M-1\}$   $M$ -ary.  $\underline{U} = (U_1, \dots, U_L)$ ,  $\underline{X} = (X_1, \dots, X_N)$ ,  $\underline{Y} = (Y_1, \dots, Y_N)$ ,  $\underline{V} = (V_1, \dots, V_N)$ . Source  $\{U_\ell\}$  is stationary with entropy rate  $H$ . Channel is DMC with  $Q = [Q(y|x)]$ .

• Define

•  $P_{e,\ell} = \Pr[V_\ell \neq U_\ell]$ .

• Average error probability/frequency:  $\bar{P}_e = \frac{1}{L} \sum_{\ell=1}^L P_{e,\ell}$  = expected frequency of errors.

•  $C = \max_{\{p(x)\}} I(p, Q)$  = **channel capacity in bits per channel use.**

•  $C' = \frac{N}{L} C$  = **channel capacity in bits per source letter.**

• Information transmission theorem for stationary source and discrete memoryless channel

1) If  $H > C' = \frac{n}{L} C$  [bit per source letter], then  $\bar{P}_e > 0$ .

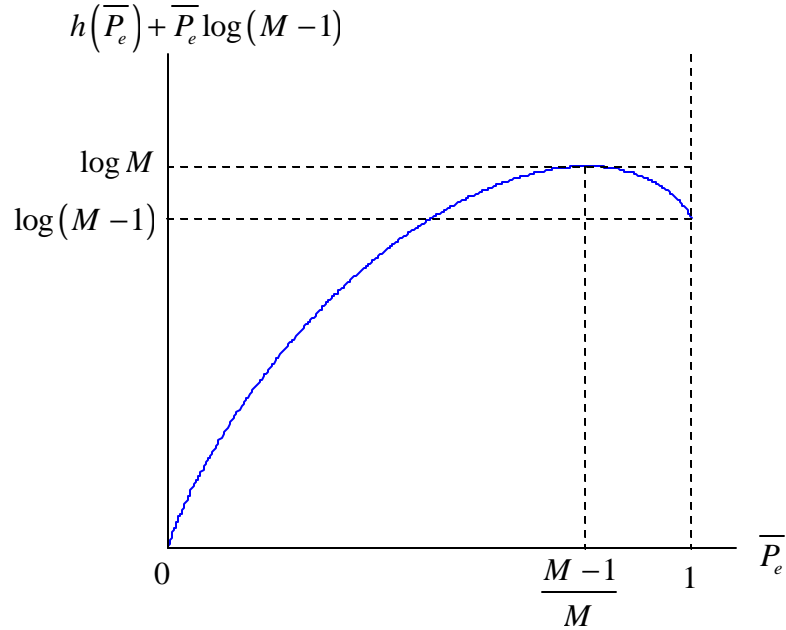
2) For any  $R < C$  and any  $\epsilon > 0$ , we can find a code  $(M = 2^{nR}, n)$  of rate  $R$  and sufficiently large block length  $n$  for which  $\max_j P_e(j) < \epsilon$ .

### Weak Converse info transmission theorem when channel is DMC

• Let  $U_1^L, V_1^L \in \mathcal{U}^L = \mathcal{V}^L$ .

•  $H - \frac{nC}{L} = H - C' \leq \frac{H(U_1^L | V_1^L)}{L} \leq h(\bar{P}_e) + \bar{P}_e \log(M-1) \leq h(\bar{P}_e) + \bar{P}_e n R_c$

• If  $H > C'$ , then  $\bar{P}_e > 0$ .



- Information transmission with an arbitrary small expected frequency of errors is not possible if the source entropy  $H = \lim_{L \rightarrow \infty} \frac{H(U_1^L)}{L}$  [bit per source symbol] exceeds the channel capacity  $C'$  measured in bit per source symbol. This conclusion holds even if one permits unbounded computation effort and is willing to tolerate enormous coding delay ( $L \rightarrow \infty$  and  $N \rightarrow \infty$  with  $L/N$  kept fixed).

### Typical Sequences for i.i.d. $X_i$

- Let  $\mathcal{X}$  be a discrete alphabet.  $\{p(x), x \in \mathcal{X}\}$  be a p.m.f.  $\bar{X} = (X_1, \dots, X_n) = X_1^n$ .  $\bar{x} = (x_1, \dots, x_n) = x_1^n \in \mathcal{X}^n$ .
- $N(x|\bar{x}) = |\{k : 1 \leq k \leq n, x_k = x\}|$ .
  - $\forall \bar{x} \in \mathcal{X}^n, \sum_{x \in \mathcal{X}} N(x|\bar{x}) = n$ .
- The **composition** / **type** of  $\bar{x} = \{N(x|\bar{x}) : \forall x \in \mathcal{X}\}$ .
- Def. Given  $\mathbf{d}$ ,  $\bar{x}$  is **d-typical** of  $\{p(x), x \in \mathcal{X}\}$  if  $\forall x \in \mathcal{X}, \left| \frac{N(x|\bar{x})}{n} - p(x) \right| < \frac{\mathbf{d}}{|\mathcal{X}|}$ .
- $T_{\mathbf{d}} = T_{\mathbf{d}}(p) = \{\bar{x} \in \mathcal{X}^n : \bar{x} \text{ is } \mathbf{d}\text{-typical}\}$ .
- $p_{\bar{x}}(\bar{x}) = \prod_{k=1}^n p(x_k) = \prod_{x \in \mathcal{X}} p(x)^{N(x|\bar{x})}$ . Thus, word probabilities  $p_{\bar{x}}(\bar{x})$  depend only on word type.
- $\bar{x} \in T_{\mathbf{d}}(p) \Rightarrow$ 
  - $\forall x \in \mathcal{X}, \left| \frac{N(x|\bar{x})}{n} - p(x) \right| < \frac{\mathbf{d}}{|\mathcal{X}|}$

- $\left| -\frac{\log p_{\bar{x}}(\bar{x})}{n} - H(\{p(x)\}) \right| < \mathbf{d} |\log p_{\min}| = \mathbf{e}_d.$ 
  - $\mathbf{e}_d > 0.$
  - $\mathbf{e}_d$  can be made arbitrary small by making  $\mathbf{d}$  small enough.  $(\lim_{d \rightarrow 0} \mathbf{e}_d = 0).$
- Define  $H^+ = H(\{p(x)\}) + \mathbf{e}_d, H^- = H(\{p(x)\}) - \mathbf{e}_d.$
- $2^{-nH^+} < p_{\bar{x}}(\bar{x}) < 2^{-nH^-}.$
- $|T_d(p)| \leq 2^{nH^+}$
- Weak Law of Large Number:  $\forall x \in \mathcal{X} \forall \mathbf{e} > 0 \lim_{n \rightarrow \infty} \Pr \left[ \left| \frac{N(x|\bar{X})}{n} - p(x) \right| > \mathbf{e} \right] = 0.$
- $\forall \mathbf{d} > 0 \lim_{n \rightarrow \infty} \Pr[\bar{X} \in T_d] = 1. \lim_{n \rightarrow \infty} \Pr[\bar{X} \notin T_d] = 0.$

### Jointly Typical Sequences for i.i.d. $(X_i, Y_i)$

- Let  $\{P(x, y) = p(x)Q(y|x), x \in \mathcal{X}, y \in \mathcal{Y}\}$  be the joint pmf over  $\mathcal{X} \times \mathcal{Y}. |\mathcal{X} \times \mathcal{Y}| = |\mathcal{X}| |\mathcal{Y}|.$
- $N(x, y | \bar{x}, \bar{y}) = |\{k : 1 \leq k \leq n, (x_k, y_k) = (x, y)\}|$ 
  - $N(x|\bar{x}) = \sum_{y \in \mathcal{Y}} N(x, y | \bar{x}, \bar{y})$
- $(\bar{x}, \bar{y})$  are **jointly  $\mathbf{d}$ -typical** of  $\{P(x, y) = p(x)Q(y|x), x \in \mathcal{X}, y \in \mathcal{Y}\}$  iff

$$\boxed{\forall x \in \mathcal{X}, \forall y \in \mathcal{Y} \left| \frac{N(x, y | \bar{x}, \bar{y})}{n} - P(x, y) \right| < \frac{\mathbf{d}}{|\mathcal{X}| |\mathcal{Y}|}.$$

- $T_d = T_d(pQ) = \{(\bar{x}, \bar{y}) : (\bar{x}, \bar{y}) \text{ is } \mathbf{d}\text{-typical of } \{P(x, y) = p(x)Q(y|x)\}\}$
- $\lim_{n \rightarrow \infty} \Pr[(\bar{X}, \bar{Y}) \notin T_d] = 0$
- $(\bar{x}, \bar{y}) \in T_d(pQ) \Rightarrow$ 
  - $2^{-n(H(P)+\mathbf{e}_d)} < P(\bar{x}, \bar{y}) < 2^{-n(H(P)-\mathbf{e}_d)}$  where  $\mathbf{e}_d = \mathbf{d} |\log P_{\min}(x, y)|.$
  - $\bar{x} \in T_d(p), \bar{y} \in T_d(q) : \text{Jointly Typicality} \Rightarrow \text{Marginal Typicality}.$

If  $(\bar{x}, \bar{y})$  is  $\delta$ -typical of  $\{P(x, y)\}$ , then

$$\bar{x} \text{ is } \delta\text{-typical of } \left\{ p(x) = \sum_y P(x, y) \right\}, \text{ and } \bar{y} \text{ is } \delta\text{-typical of } \left\{ q(y) = \sum_x P(x, y) \right\}.$$

- $\Rightarrow p(\bar{x}) \leq 2^{-nH(p)^-}$ ,  $q(\bar{y}) \leq 2^{-nH(q)^-}$
- $|T_d(P)| \leq 2^{nH(P)^+}$ .

### Direct Noisy-Channel Coding Theorem for d.m.c.

- **Block code**:  $\mathcal{C} = \{\bar{x}^{(1)}, \dots, \bar{x}^{(M)}\}$ .  $\bar{x}^{(i)} \in \mathcal{X}^n$ . **Rate of the code** =  $R = \frac{\log M}{n}$ .  $M = 2^{nR}$ . Hence,  $(2^{nR}, n)$  code.

- **Random code selection / random coding**: generate  $\mathcal{C}$  at random according to the distribution  $p = \{p(x), x \in \mathcal{X}\}$ :

Let  $\mathcal{C} = \{\bar{X}^{(1)}, \dots, \bar{X}^{(M)}\}$  be a randomly chosen block code such that all  $nM$  letters (of the  $M$  codewords each with  $n$  letters) are i.i.d.  $\{p(x)\}$ .

- Note that once code is select, you use it in a deterministic way.
- $\Pr[\bar{X}^{(j)} = x_1^n] = \prod_{k=1}^n p(x_k)$
- $\Pr[X_1^n = x_1^n] = \prod_{k=1}^n p(x_k)$
- Thus, the channel input is i.i.d. with  $p = \{p(x), x \in \mathcal{X}\}$ .
- $\Pr[\bar{X}^{(j)} = \bar{x} | J = j] = \Pr[\bar{X}^{(j)} = \bar{x}] = \prod_{k=1}^n p(x_k)$
- Independent of the random code, let  $J$  be the **random message index** with pmf  $\{p_j, 1 \leq j \leq M\}$ . If  $J = j$ , then the components  $\bar{X}_1^{(j)}, \bar{X}_2^{(j)}, \dots, \bar{X}_n^{(j)}$  of  $\bar{X}^{(j)}$  will be put into the channel in this order during  $n$  successive channel uses.
  - So, by knowing  $J$ , we know  $\log(M)$  bits in  $n$  channel uses.
  - $\bar{Y} = (Y_1, \dots, Y_n) \equiv$  the resulting channel output vector.
  - $\hat{J} \equiv \hat{J}(\bar{Y}) =$  the decoder's estimate of  $J$  based on  $\bar{Y}$ .
  - $\bar{P}_e = \Pr[\hat{J} \neq J]$  which depends on the joint distribution of  $J$ ,  $\mathcal{C}$ ,  $\bar{Y}$ , and  $D$ .

• **Joint typicality decoding rule**: Upon receiving  $\bar{y}$ , if  $\exists! j^* 1 \leq j^* \leq M$   $(\bar{x}_{j^*}, \bar{y}) \in T_d$ , decode  $\bar{y}$  by  $\hat{J}(\bar{y}) = j^*$ . If no such index or if there is more than one such index, declare a decoding error.

- $\bar{P}_e(j) = \Pr[\hat{J} \neq J | J = j]$  (averaged over all code).
  - By symmetry of code selection scheme,  $\forall j \bar{P}_e(j) = \bar{P}_e(1)$ .
- Overall average error probability:  $\bar{P}_e = \Pr[\hat{J} \neq J] = \sum_{j=1}^M p_j \bar{P}_e(j) = \bar{P}_e(1)$ .



- Probability facts:

- $(\bar{x}, \bar{y}) \in T_d(pQ) \Rightarrow p_{\bar{X}^{(\ell)}}(\bar{x}) \leq 2^{-nH(p)^-}$ , and  $q_{\bar{Y}}(\bar{y}) \leq 2^{-nH(q)^-}$

- For  $\ell = 1$

- $Q_{\bar{Y}|\bar{X}^{(1)}, J}(\bar{y}|\bar{x}, 1) = \prod_{k=1}^n Q(y_k|x_k)$ .

- $P_{\bar{X}^{(1)}, \bar{Y}|J}(\bar{x}, \bar{y}|1) = \prod_{k=1}^n p(x_k)Q(y_k|x_k)$ .

- For  $\ell \neq 1$ ,

- $\{\bar{X}^{(\ell)}, 2 \leq \ell \leq M\}$  is independent of  $\bar{Y}$ .

- $Q_{\bar{Y}|\bar{X}^{(\ell)}, J}(\bar{y}|\bar{x}, 1) = Q_{\bar{Y}}(\bar{y}) = \prod_{k=1}^n q(y_k)$

- $P_{\bar{X}^{(\ell)}, \bar{Y}|J}(\bar{x}, \bar{y}|1) = p_{\bar{X}^{(\ell)}}(\bar{x})q_{\bar{Y}}(\bar{y}) = p_{\bar{X}}(\bar{x})q_{\bar{Y}}(\bar{y})$  same for all  $\ell \neq 1$ .

- $\Pr\{(\bar{X}^{(\ell)}, \bar{Y}) \in T_d(P)|J=1\} = \Pr\{(\bar{X}^{(2)}, \bar{Y}) \in T_d(P)|J=1\} \leq 2^{nH(pQ)^+} 2^{-nH(p)^-} 2^{-nH(q)^-} = 2^{-nI(p, Q)^-}$ , where

$$I(p, Q)^- = I(p, Q) - \mathbf{e}_d \text{ where } \mathbf{e}_d = \mathbf{d} \left( \left| \log \min_x p(x) \right| + \left| \log \min_y q(y) \right| + \left| \log \min_{x,y} p(x)Q(y|x) \right| \right).$$

- $\bar{P}_e(1) = \Pr[E_1 \cup E_2] \leq \Pr[E_1] + \Pr[E_2]$

- $E_1 = \{(\bar{X}_1, \bar{Y}) \notin T_d(pQ)|J=1\}$ .  $\lim_{n \rightarrow \infty} \Pr[E_1] = 0$  because  $(\bar{X}_1, \bar{Y}) \sim pQ$ .

- $E_2 = \bigcup_{\ell > 1} \{(\bar{X}_\ell, \bar{Y}) \in T_d(pQ)|J=1\}$ .  $\Pr[E_2] \leq 2^{-n(I(p, Q)^- - R)}$ .

- If  $I(p, Q)^- > R$ ,  $\lim_{n \rightarrow \infty} \Pr[E_2] = 0$ .

- The channel coding theorem:

- Direct: All rates below capacity  $C$  are achievable.

- $\forall R, R < C$ , there exists a sequence of  $(2^{nR}, n)$  codes along which the error probability decays to zero as  $n \rightarrow \infty$ . (regardless of the probabilities  $P_j$  of the messages.)

- $\forall R, R < C, \forall \mathbf{e} > 0 \exists$  a  $(2^{nR}, n)$  code of rate  $R$  and sufficiently large block length  $n$  for which  $\max_j P_e(j) < \mathbf{e}$ .

- Conversely, any sequence of  $(2^{nR}, n)$  codes with  $\lim_{n \rightarrow \infty} \mathbf{I}^n = 0$  must have  $R \leq C$ .

**Etc**

- Separation Theorem for source and channel coding.

Let  $\{U_k\}$  be an ergodic stationary information source with entropy rate  $H$ . If  $H < C$ , the capacity of the DMC  $Q(y|x)$ , then it is possible to convey  $\underline{U}$  through the channel with an arbitrary small probability of error. Employ a source code with rate  $R_s$  and a channel code with rate  $R_c$  such that  $H < R_s < R_c < C$  [bit/sec].

- Asymptotic optimality can be achieved by separating source coding and channel coding.
- The source code and the channel code can be designed separately without losing asymptotic optimality