

Tuesday, November 27th, 2012

Chapter 6: Polar codes

I - Motivation

We will restrict our discussion to channels with binary inputs $x \in \{0, 1\}$ and y can be anything

Among these channels there are two "extremal" channels over which it is easy to communicate at capacity:

1) Completely useless channels where $C=0$. These channels have $p(y|0) = p(y|1)$ (output independent of input).

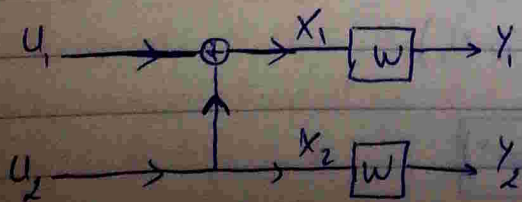
2) Completely noiseless channels where $C=1$. These channels have $p(y|0)p(y|1) = 0$ (output determines the input).

Polar codes try to turn a given channel to one of the above "extremal" channels where it is easy to communicate at capacity.

II - Polar transform

Polar transform is a method to create extremal channels from multiple-uses of a given channel W .

1. 2×2 building block



Set $X_1 = U_1 \oplus U_2$ (mod 2 sum)
 $X_2 = U_2$

Suppose U_1, U_2 are independent and uniformly distributed on $\{0, 1\}$

$\Rightarrow (X_1, X_2)$ are also independent, uniformly distributed on $\{0, 1\}^2$

$$\begin{aligned} I(U_1, U_2; Y_1, Y_2) &= I(X_1, X_2; Y_1, Y_2) \quad \text{since } (U_1, U_2) \text{ are in one-to-one} \\ &\quad \text{correspondence with } (X_1, X_2) \\ &= I(X_1; Y_1) + I(X_2; Y_2) \quad \text{since channel is memoryless} \\ &\quad \text{and input independent} \\ &= 2I(w) \end{aligned}$$

where $I(w)$ is the mutual information between input and output of w when input is uniform on $\{0, 1\}$

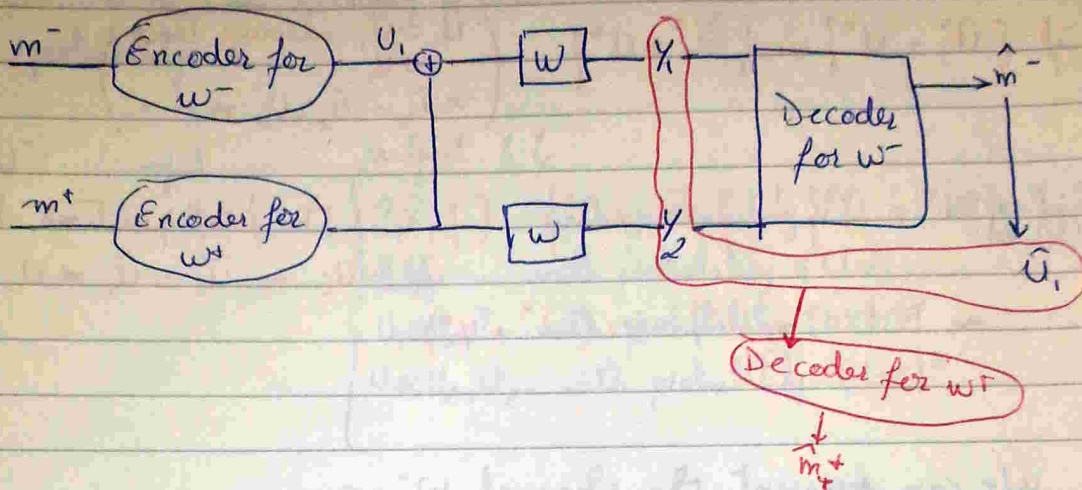
$$\begin{aligned} \text{So far } 2I(w) &= I(U_1, U_2; Y_1, Y_2) \\ &= I(U_1; Y_1, Y_2) + I(U_2; Y_1, Y_2 / U_1) \quad \text{chain rule} \\ &= I(U_1; Y_1, Y_2) + I(U_2; Y_1, Y_2 / U_1) \quad \text{since } U_1, U_2 \\ &\quad \text{independent} \\ &\Rightarrow I(U_2; Y_1, Y_2 / U_1) \\ &= I(U_2; U_1) + I(U_2; Y_1, Y_2 / U_1) \end{aligned}$$

Consider the two synthetic channels:

- 1/ w^- : input U_1 , output (Y_1, Y_2)
- 2/ w^+ : input U_2 , output (Y_1, Y_2, U_1)

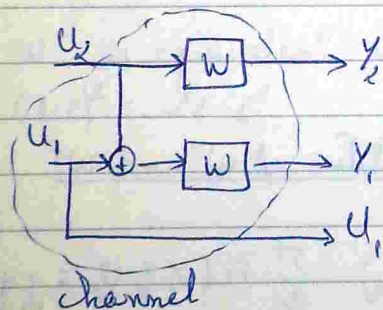
$$\Rightarrow 2I(w) = I(w^-) + I(w^+)$$

However to obtain U_1 as an output of W^+ we can impose a successive decoder order at the receiver



See after 2 pages

So we can represent channel W^+ as:



$$I(W^+) = I(U_2, Y_2 | U_1) \geq I(U_2, Y_2) = I(W)$$

$$\Rightarrow I(W) \leq I(W) \leq I(W^+)$$

Remark: On "helped" and "unhelped" decoding

Helped decoding:

We decode U_1, U_2, \dots, U_m from some observation Y_1, Y_2, \dots, Y_n

$$\hat{U}_1 = \phi_1(Y_1 \rightarrow Y_n)$$

$$\hat{U}_2 = \phi_2(Y_1 \rightarrow Y_n, \hat{U}_1)$$

$$\hat{U}_i = \phi_i(Y_1 \rightarrow Y_n, \hat{U}_1 \rightarrow \hat{U}_{i-1})$$

$$\hat{U}_m = \phi_m(Y_1 \rightarrow Y_n, \hat{U}_1 \rightarrow \hat{U}_{m-1})$$

Unhelped decoding

We decode $U_1 \rightarrow U_m$ from some observation $Y_1 \rightarrow Y_n$

$$\tilde{U}_1 = \phi_1(Y_1 \rightarrow Y_n) = \hat{U}_1$$

$$\tilde{U}_2 = \phi_2(Y_1 \rightarrow Y_n, \tilde{U}_1)$$

$$\tilde{U}_i = \phi_i(Y_1 \rightarrow Y_n, \tilde{U}_1 \rightarrow \tilde{U}_{i-1})$$

$$\tilde{U}_m = \phi_m(Y_1 \rightarrow Y_n, \tilde{U}_1 \rightarrow \tilde{U}_{m-1})$$

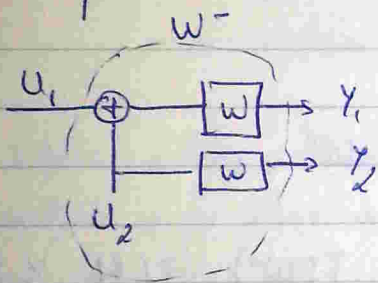
$P_2(\hat{u}^n \neq u^n)$ = probability of error of helped system

$P_2(\tilde{u}^n \neq u^n)$ = probability of error of unhelped system

$$\neg \{ \hat{u}^n = u^n \} \Rightarrow \{ \tilde{u}^n \neq u^n \} \quad \left(\begin{array}{l} \hat{u}_i = \hat{u}_i \text{ always true} \\ = u_i \\ \Rightarrow \tilde{u}_2 = \hat{u}_2 = u_2 \text{ etc } \dots \end{array} \right)$$

$$\neg \{ \hat{u}^n \neq u^n \} \Rightarrow \{ \tilde{u}^n \neq u^n \} \quad (\text{if } i \text{ is the first position where } \hat{u}_i \neq u_i \Rightarrow \tilde{u}_i = \hat{u}_i \neq u_i)$$
$$\Rightarrow P_2(\hat{u}^n \neq u^n) = P_2(\tilde{u}^n \neq u^n)$$

We can represent the channel w^- as:



Channel transition probabilities of w^- and w^+

$$w^-(y_1, y_2 / u_1) = P_2(y_1, y_2 = y_1, y_2 / u_1 = u_1) = \sum_{u_2} P_2(y_1, y_2, u_2 = y_1, y_2, u_2 / u_1 = u_1)$$

$$= \sum_{u_2} \underbrace{P_2(u_2 = u_2 / u_1 = u_1)}_{= 1/2} P_2(y_1, y_2 = y_1, y_2 / u_1, u_2 = u_1, u_2)$$

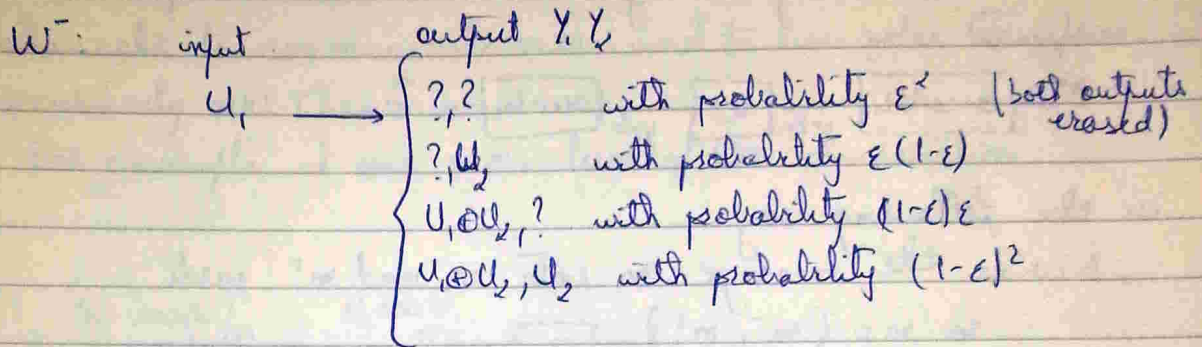
$$= \frac{1}{2} \sum_{u_2} w(y_1 / u_1 \oplus u_2) w(y_2 / u_2)$$

$$= \frac{1}{2} w(y_1 / u_1) w(y_2 / 0) + \frac{1}{2} w(y_1 / u_1 \oplus 1) w(y_2 / 1)$$

Similarly $w^+(y_1, y_2, u_1 / u_2) = \frac{1}{2} w(y_1 / u_1 \oplus u_2) w(y_2 / u_2)$

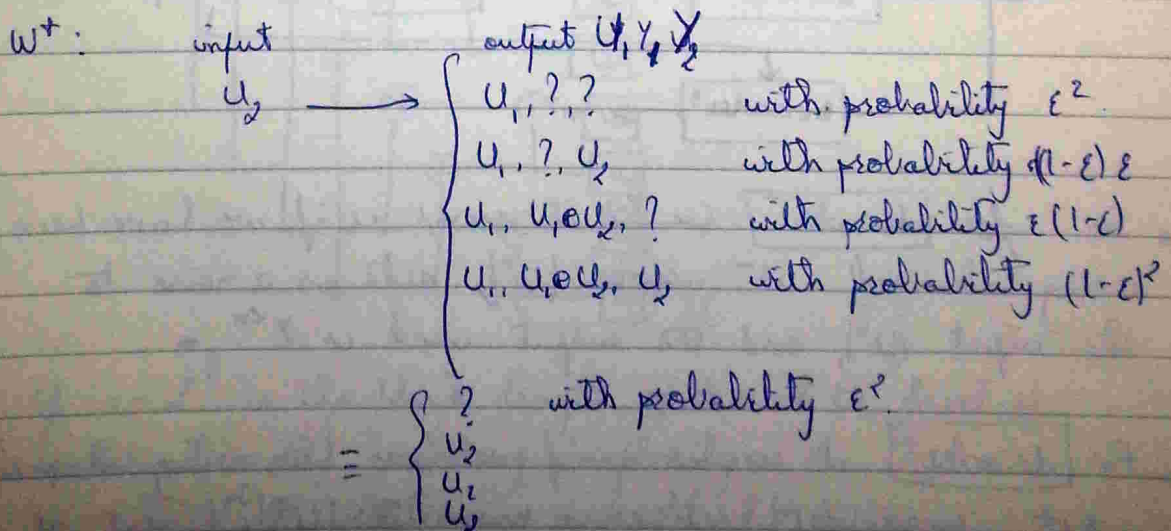
2- Example 1

Let W be the Binary Erasure Channel with erasure probability ϵ



$$= \begin{cases} ? \\ ? \\ ? \\ u_1 \end{cases} \text{ with probability } (1-\epsilon)^2$$

So W^- is equivalent to BEC with erasure probability $1 - (1-\epsilon)^2 = 2\epsilon - \epsilon^2$

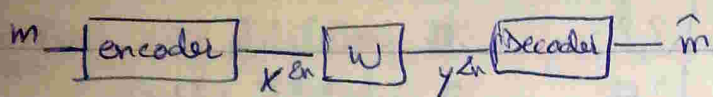


So W^+ is equivalent to BEC with erasure probability ϵ^2

3- Explanation of successive decoder

We want to develop a way to implement channels w^- and w^+ .

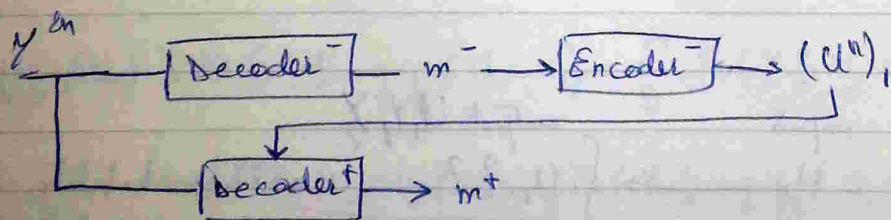
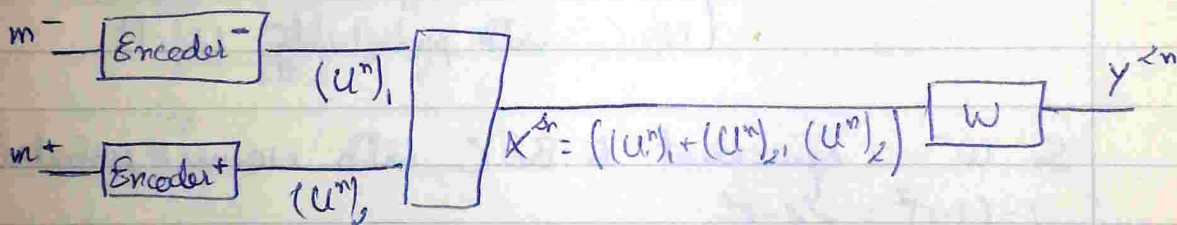
So we want to design a system



in the following way:

- Divide the message m into m^- and m^+ parts so $m = (m^-, m^+)$

- Consider the following construction



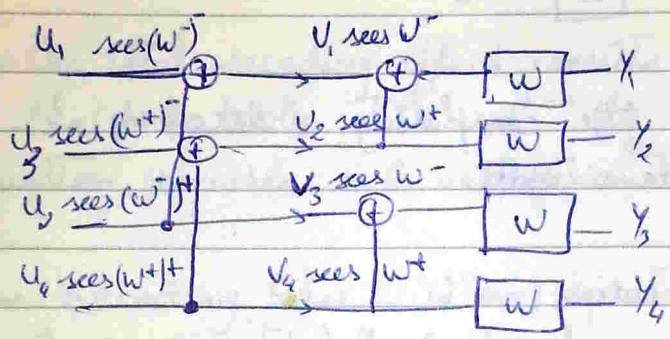
- So for decoder^- (we had) it seems as if we have been using channel w^- since $(U^n)_2$ acts as a noise to the input $(U^n)_1$, and the output used is y^n .
- For decoder^+ it seems as if we have been using channel w^+ since $(U^n)_1$ acts as a noise to $(U^n)_2$ (input of the channel). Moreover decoder^+ uses both y^n

and (u^n) generated at receiver to decode m^+ . So it is as if the output of the channel is $(y^{2n}, (u^n)_1)$

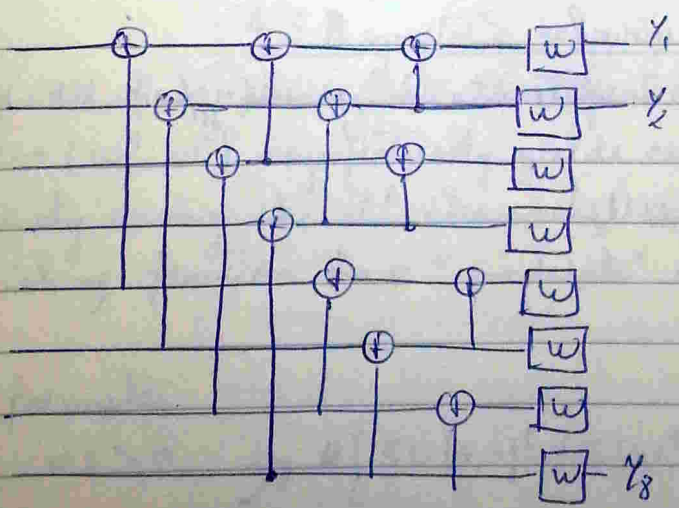
4. Recursive Channels

The hope is that the channels w^- and w^+ will be easier channels to communicate over, and to use polar transforms recursively over them.

For example at one more level of the construction



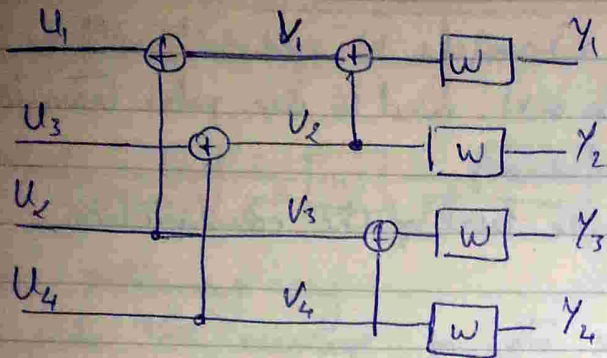
4.1. Encoding capacity



if $n = 2^l$ is the block length of the l 'th level of polarization, the encoding complexity is:
 $\frac{2^l}{2} l$ additions and $\frac{2^l}{2} l$ copies
 so we have $O(n \log n)$ operations

So we see that encoding is cheap.

4.2. Decoding complexity



In order to calculate the complexity of decoding we need to follow a certain order.

- First we start by deciding on u_1 .

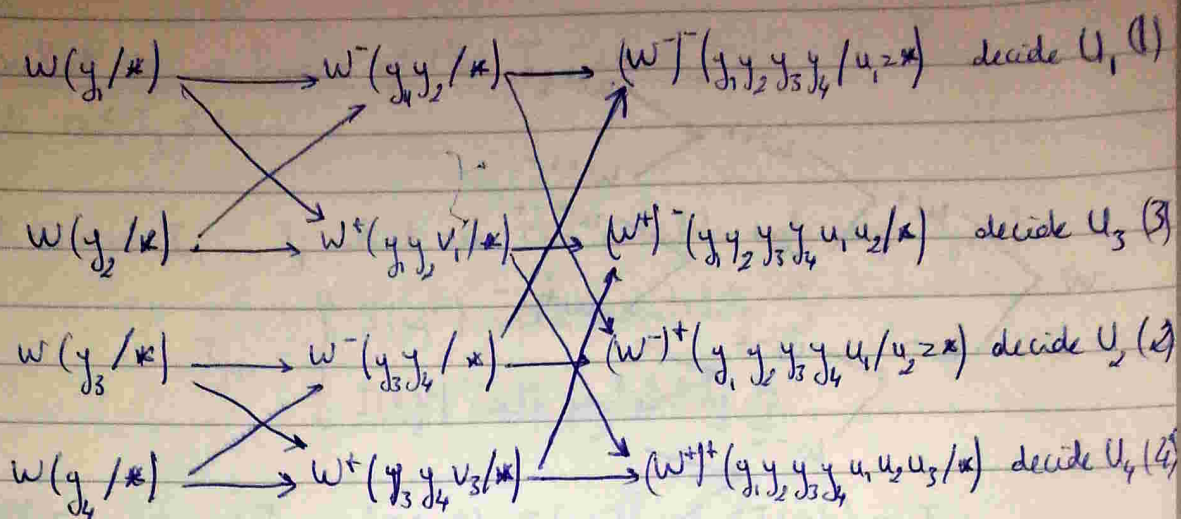
To decide on u_1 , we need to find $(w^-)(y_1, y_2, y_3, y_4 / u_1 = *)$. This quantity is determined by $w^-(y_1, y_2 / *)$ and $w^-(y_3, y_4 / *)$.

Similarly $w^-(y_1, y_2 / *)$ is determined by $w^-(y_1 / *)$ and $w^-(y_2 / *)$.

- Second we decide on u_2 .

- we decide on u_3 .

- We decide u_4 .



So the total decoding effect consists of ~~filling~~ filling in a $2^l \times (l+1)$ table with $W^{+...+}(y_1 \dots y_l/x)$, each filling equation requires ≤ 6 arithmetic equations

\Rightarrow Decoding takes $O(n \log n)$ equations
So decoding is also cheap.

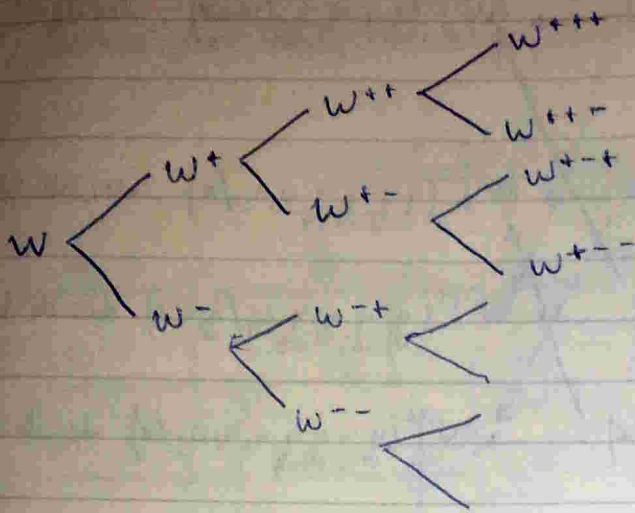
4.3. Asymptotic behavior of recursive channels

We will now show that the repeated application of the $W \rightarrow (W^+, W^-)$ asymptotically yields extremal channels. Namely among the 2^l channels $W^{+...+}, \dots, W^{-...-}$ a vanishing fraction have "moderate" value of $I(W)$.

In formula,

$$\forall \epsilon > 0, \frac{1}{2^l} \#\{W^{\pm \dots \pm} \mid I(W^{\pm \dots \pm}) \in (s, 1-s)\} \rightarrow 0 \text{ as } l \text{ gets large}$$

For this, let us first organize the W^{\pm} 's in a binary tree:



Consider climbing the tree randomly:

$$w_0 = w$$

$$w_{i+1} = \begin{cases} w_i^+ & \text{with probability } \frac{1}{2} \\ w_i^- & \text{with probability } \frac{1}{2} \end{cases}$$

w_l is a randomly chosen channel among the 2^l channels at level l of the tree.

So $I_l = I(w_l)$ is a random variable taking values in $[0, 1]$

$$\text{and } P_z(I(w_l) \in (s, 1-s)) = \frac{1}{2^l} \# \{ \bar{s} \in \{+, -\}^l / I(w^{\bar{s}}) \in (s, 1-s) \}$$

Example: $I_0 = I(w)$, $I_1(w) = \begin{cases} I(w^+) & \text{with probability } \frac{1}{2} \\ I(w^-) & \text{with probability } \frac{1}{2} \end{cases}$

$$I_2(w) = \begin{cases} I(w^{++}) & \text{with probability } \frac{1}{4} \\ \vdots \\ I(w^{--}) & \text{with probability } \frac{1}{4} \end{cases}$$

What do we know of the I_j process?

1) $0 \leq I_j \leq 1$

2) Given I_0, \dots, I_j what are the possible values of I_{j+1} ?

$$I_{j+1} = \begin{cases} I(w_j^+) & \text{with probability } \frac{1}{2} \\ I(w_j^-) & \text{with probability } \frac{1}{2} \end{cases}$$

$$E[I_{j+1} | I_0, \dots, I_j] = \frac{1}{2} [I(w_j^+) + I(w_j^-)] = \frac{1}{2} \cdot 2 I(w_j) = I_j$$

From the homework we know that for a real process

$$E[(I_{k+1} - I_k)(I_{j+1} - I_j)] = 0 \quad \text{when } k \neq j$$

Consequently $1 \geq |I_l - I_0|^2$ because both I_0 and I_l are in $[0, 1]$

$$\begin{aligned} \Rightarrow E(1) &\geq E[|I_l - I_0|^2] \\ &= E\left[\left(\sum_{j=0}^{l-1} I_{j+1} - I_j\right)^2\right] \\ &= \sum_{j=0}^{l-1} \underbrace{E[(I_{j+1} - I_j)^2]}_{\geq 0} + \sum_{j \neq k} \cancel{E[(I_{j+1} - I_j)(I_{k+1} - I_k)]} \end{aligned}$$

$$\Rightarrow 1 \geq \sum_{j=0}^{l-1} E[(I_{j+1} - I_j)^2] \quad \text{for any } l.$$

$$\text{Let } \sigma_j^2 = E[(I_{j+1} - I_j)^2]$$

$$\Rightarrow \sum_{j=0}^{l-1} \sigma_j^2 \leq 1 \quad \text{for every } l$$

$$\Rightarrow \lim_{l \rightarrow \infty} \sigma_l^2 = 0$$

$$\Rightarrow E[|I_{l+1} - I_l|^2] \rightarrow 0$$

$$\Rightarrow \Pr(|I_{eH} - I_e| > \varepsilon)$$

$$\leq \frac{E(|I_{eH} - I_e|^2)}{\varepsilon^2}$$

\Rightarrow for any $\varepsilon > 0$, $\Pr(|I_{eH} - I_e| > \varepsilon) \rightarrow 0$ as l gets large

this is equivalent to $\frac{1}{2^l} \#\{\bar{S} \in \mathcal{S}_+^l / I(w^{\bar{S}}) - I(w^{\bar{S}^-}) > \varepsilon\} \rightarrow 0$
as l gets large

Interpretation: At level l the reset channel is very close to the present one in terms of mutual information

Lemma 1

$\forall \delta > 0, \exists \varepsilon > 0$ such that

$$I(w) - I(w^-) \leq \varepsilon$$

$$\Rightarrow I(w) \notin (\delta, 1-\delta) \quad (\text{Proof after 3 pages})$$

This means that if the channel does not "move by much" then it must be an extremal channel (either good or bad).

Using this lemma 1 we can deduce that

$$\frac{1}{2^l} \#\{\bar{S} \in \mathcal{S}_+^l / I(w^{\bar{S}}) \in (\delta, 1-\delta)\} \rightarrow 0 \text{ as } l \text{ gets large}$$

\Rightarrow Polarization happens

So as l gets large the channel becomes extremal and hence this polarization

Further note:

$$I(w^-) + I(w^+) = 2I(w)$$
$$\underbrace{I(w^-) + I(w^{-+})}_{2I(w)} + \underbrace{I(w^{+-}) + I(w^+)}_{2I(w)} = 4I(w)$$

$$\Rightarrow \frac{1}{2^l} \sum_{s \in \{+, -\}^l} I(w^s) = I(w) \quad \text{for every } l.$$

Theorem 1:

For any $s > 0$,

$$(1) \frac{1}{2^l} \#\{s \in \{+, -\}^l / I(w^s) \in [1-s, 1]\} \rightarrow I(w) \quad (\text{Number of good channels})$$

$$(2) \frac{1}{2^l} \#\{s \in \{+, -\}^l / I(w^s) \in [0, s]\} \rightarrow 1 - I(w) \quad (\text{Number of bad channels})$$

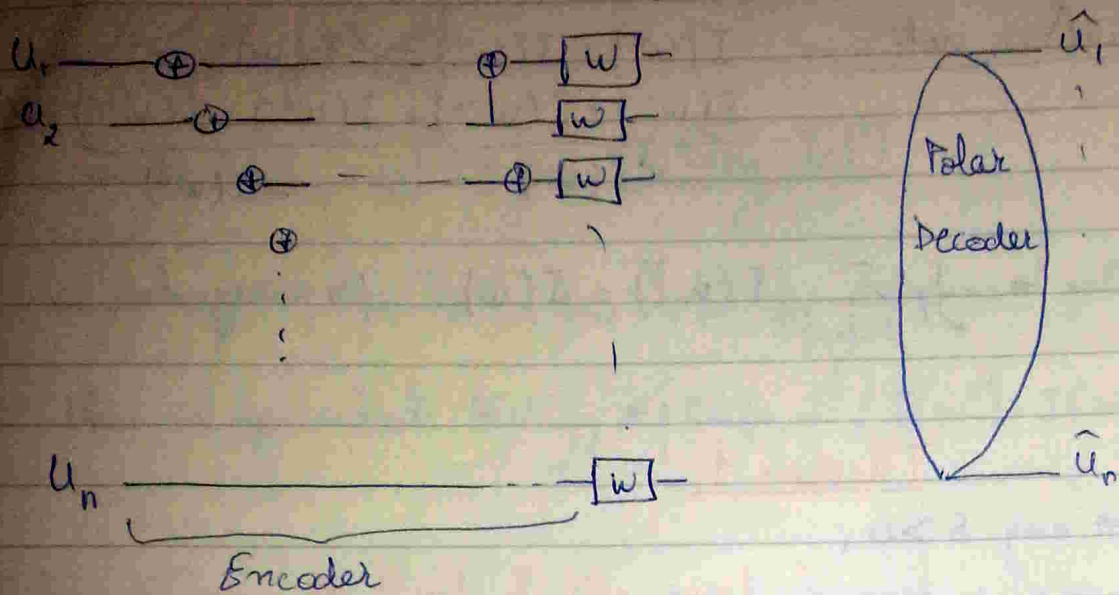
$$(3) \frac{1}{2^l} \#\{s \in \{+, -\}^l / I(w^s) \in (s, 1-s)\} \rightarrow 0 \quad (\text{Number of moderate channels})$$

4.4. Alternative design for communication systems

The results of the previous section suggest a method of designing a communication system to transmit data over w :

Given $R < I(w)$, $s > 0$, pick l large enough so that $\frac{1}{2^l} \#\{s \in \{+, -\}^l / I(w^s) \geq 1-s\} \geq R$

Then construct a polar encoder/decoder for polarization depth l , with $n = 2^l$



- Once this is done, we will have i 's such that u_i sees a good channel (there are $\geq R \cdot n$ of them)
Use these indices to send data.
- The remaining u_i 's freeze them to randomly chosen values in $\{0, 1\}$ known to both sender and receiver.

This gives us a code that operates at rate $\geq R$ and where each data bit travels on a synthetic channel with mutual information $\geq 1 - \epsilon$.

N.B.: The input to bad channels are chosen randomly rather than in a deterministic way since in our calculations of $I(w^-)$ and $I(w^+)$ we assumed the inputs u_1 and u_2 to be iid drawn from $\{0, 1\}$ according to uniform distribution

Remarks

✓ For "symmetric channels" there is no need to choose the frozen u_i 's by a random experiment, they can be chosen to equal 0.

Definition 1

A binary input channel W is said to be symmetric if there is a function

$$\pi: \mathcal{Y} \rightarrow \mathcal{Y} \text{ with}$$

$$(1) \pi(\pi(y)) = y \quad (\pi \text{ is its own inverse})$$

$$(2) W(y|0) = W(\pi(y)|1)$$

Example:

• BSC is symmetric $(\pi: 0 \rightarrow 1, 1 \rightarrow 0)$

• BEC is symmetric $(\pi: 0 \rightarrow 1, 1 \rightarrow 0, \mathcal{E} \rightarrow \mathcal{E})$

• Binary input additive Gaussian channel

$$\mathcal{X} = \{0, 1\}, \quad \mathcal{Y} = \mathbb{R}$$

$$Y = (-1)^X + Z \quad Z \sim \mathcal{N}(0, \sigma^2)$$

Check that ~~$\pi(y) = -y$~~ $\pi(y) = -y$ satisfies (1) and (2) in the symmetry definition, i.e. this channel is also symmetric.

III. On the error probability of polar codes

Recall the difference between decoding with help and the unhelped decoding

$$\hat{u}_1 = \phi_1(y_1 \dots y_n)$$

$$\hat{u}_2 = \phi_2(y_1 \dots y_n, u_1)$$

⋮

$$\hat{u}_n = \phi_n(y_1 \dots y_n, u_1, \dots, u_{n-1})$$

so $P_z(\hat{u}_i \neq u_i) =$ probability of error if a bit is sent on $w^{\bar{s}}$
(where \bar{s} is the sequence of $\{+, -\}$ corresponding to u_i)

Example: $n=8$, u_1 travels on w^{---} , $u_1 \rightarrow y^8$
 u_2 travels on w^{--+} , $u_2 \rightarrow y^8 u_1$
 u_3 travels on w^{-+-} , $u_3 \rightarrow y^8 u_1 u_2$
⋮

The unhelped decoder (implemented by the polar decoder).

$$\tilde{u}_1 = \phi_1(y^n)$$

$$\tilde{u}_2 = \phi_2(y^n, \tilde{u}_1)$$

⋮

$$\tilde{u}_n = \phi_n(y^n, \tilde{u}^{n-1})$$

We had shown that

$$P_z(\tilde{u}^n \neq u^n) = P_z(\hat{u}^n \neq u^n) \leq \sum_{i=1}^n P_z(\hat{u}_i \neq u_i)$$

$$P_z(\hat{u}_i \neq u_i) = \begin{cases} 0 & i \text{ is frozen} \\ P_e(w^{\bar{s}}) & \bar{s} \text{ is the index of the channel on} \\ & \text{which } u_i \text{ is sent.} \end{cases}$$

$$\text{with } P_e(w) = \frac{1}{2} \sum_y w(y/0) \mathbb{1}\{w(y/1) \geq w(y/0)\} + \frac{1}{2} \sum_y w(y/1) \mathbb{1}\{w(y/0) \geq w(y/1)\}$$

It is not difficult to show that
 $I(w) \geq 1 - \delta \rightarrow P_e(w) \leq \delta$

Consequently,

$P(\tilde{U}^n \neq U^n)$ of the polar codes is upper bounded
 by:

$$\sum_{\substack{\bar{s}: \bar{s} \text{ is} \\ \text{an index} \\ \text{that is used} \\ \text{for data} \\ \text{transmission}}}^n \underbrace{(1 - I(w^{\bar{s}}))}_{\leq \delta_e}$$

It turns out ~~that~~ the polarization takes place faster than
 we proved. Indeed the following holds:

$$\frac{1}{2^l} \#\{\bar{s} \in \{1, \dots, 2^l\} / I(w^{\bar{s}}) \in (\delta_e, 1 - \delta_e)\} \rightarrow 0$$

$$\text{and } \delta_e = \frac{1}{2^{2\beta l}} \quad \beta < \frac{1}{2}$$

So by increasing the range of moderate codes we still get
 fast polarization

Consequently, ($n = 2^l$)

$$P(\tilde{U}^n \neq U^n) \leq \frac{1}{2^{2\beta l}} \rightarrow 0 \text{ as } l \rightarrow \infty$$

Corner Proof:

Proof of Lemma 1

Lemma 1 can be stated as follows:

if $I(w) - I(w')$ is small then $I(w)$ is extremal

Proof:

Lemma 2: if (X_1, Y_1) and (X_2, Y_2) are independent and if

X_1 and X_2 are \mathbb{F}_2 valued ($\{0,1\}$)

and suppose $H(X_1/Y_1) = h_2(p_1)$

$H(X_2/Y_2) = h_2(p_2)$

then $H(X_1 \oplus X_2 / Y_1, Y_2) \geq h_2(p_1 * p_2)$

where $p_1 * p_2 = p_1(1-p_2) + p_2(1-p_1)$

Before we prove Lemma 2, let's see how we can use it:

Recall $w: X \rightarrow Y$

$$I(w) = I(X, Y) \Big|_{X \sim \text{uniform}\{0,1\}} = H(X) - H(X/Y) = 1 - H(X/Y)$$

$w^-: U_1 \rightarrow Y_1, Y_2$

$$I(w^-) = H(U_1) - H(U_1 / Y_1, Y_2)$$

$$= 1 - H(X_1 \oplus X_2 / Y_1, Y_2)$$

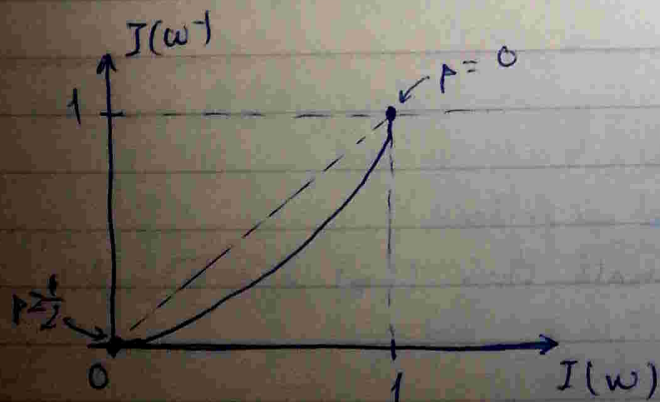
since $U_1 = X_1 \oplus X_2$ is uniform $\{0,1\}$

(X_1, Y_1) and (X_2, Y_2) are independent
identically distributed as (X, Y)

Thus if $H(X/Y) = h_2(p)$

then $H(X_1 \oplus X_2 / Y_1, Y_2) \geq h_2(p * p)$

(by lemma 2 called "Mrs Gerber's Lemma")



$$\text{if } I(w) = 1 - h_2(p)$$

$$\approx I(w^-) \leq 1 - h_2(p * p)$$

\Rightarrow for every $\delta > 0$, $\exists \epsilon > 0$ such that $I(w^-) \geq I(w) - \epsilon \approx I(w) \notin (\delta, 1 - \delta)$

This proves Lemma 1. So Mrs. Geiber's is sufficient to prove "polarization happens" ■

Proof of Mrs. Geiber's Lemma:

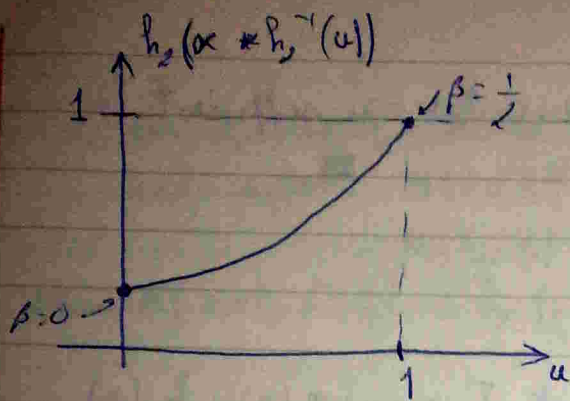
Define $\eta_1(y_1) = H(X_1 | Y_1 = y_1)$ so that $\sum_{j_1} p_{j_1}(y_1) \eta_1(y_1) = h_2(P_1)$

$\eta_2(y_2) = H(X_2 | Y_2 = y_2)$ so that $\sum_{j_2} p_{j_2}(y_2) \eta_2(y_2) = h_2(P_2)$

$$H(X_1 \oplus X_2 | Y_1 = y_1, Y_2 = y_2) = h_2 \left(\underbrace{h_2^{-1}(\eta_1(y_1))}_{P_1(X_1=1|Y_1=y_1)} * \underbrace{h_2^{-1}(\eta_2(y_2))}_{P_2(X_2=1|Y_2=y_2)} \right)$$

$$\begin{aligned} \text{so } H(X_1 \oplus X_2 | Y_1, Y_2) &= \sum_{j_1} p_{j_1}(y_1) \sum_{j_2} p_{j_2}(y_2) h_2 \left(h_2^{-1}(\eta_1(y_1)) * h_2^{-1}(\eta_2(y_2)) \right) \\ &\stackrel{(a)}{\geq} \sum_{j_1} p_{j_1}(y_1) h_2 \left(h_2^{-1}(\eta_1(y_1)) * h_2^{-1} \left(\sum_{j_2} p_{j_2}(y_2) \eta_2(y_2) \right) \right) \\ &= \sum_{j_1} p_{j_1}(y_1) h_2 \left(h_2^{-1}(\eta_1(y_1)) * h_2^{-1}(h_2(P_2)) \right) \\ &= \sum_{j_1} p_{j_1}(y_1) h_2 \left(h_2^{-1}(\eta_1(y_1)) * P_2 \right) \\ &\stackrel{(b)}{\geq} h_2 \left(h_2^{-1} \left(\underbrace{\sum_{j_1} p_{j_1}(y_1) \eta_1(y_1)}_{h_2(P_1)} * P_2 \right) \right) \\ &= h_2(P_1 * P_2) \end{aligned}$$

(a) and (b) would be resolved if the function $u \mapsto h_2(x * h_2^{-1}(u))$ is convex \cup in u because then we would have $\sum_u p(u) h_2(x * h_2^{-1}(u)) \geq h_2(x * h_2^{-1}(\sum_u p(u)u))$



Set $u = h_2(\beta)$ and plot $h_2(\beta)$ on the horizontal axis and $h_2(x * \beta)$ in the vertical axis and sweep β from 0 to $\frac{1}{2}$

The function $h_2(x * h_1^{-1}(h_2(\beta))) = h_2(x * \beta)$ is convex
 \equiv slope is increasing as we move from left to right

$$\equiv \frac{\frac{d}{d\beta} h_2(x * \beta)}{\frac{d}{d\beta} h_2(\beta)} \text{ is increasing with } \beta$$

$$\text{But } \frac{d}{d\beta} h_2(\beta) = \log \frac{1-\beta}{\beta}$$

$$\text{and } \frac{d}{d\beta} h_2(x * \beta) = (1-x\beta) \log \frac{1-(x\beta)}{x\beta}$$

and it is easy to verify that the ratio is increasing in β \square