

---

---

---

---

---



# Codage et entropie

Théorie de Shannon (1948).

Signal  $\equiv$  suite  $x_0, x_1, x_2, \dots$  d'éléments d'un ensemble fini  $A = (a_1, \dots, a_k)$  l'alphabet.

Ex: • un texte est un signal sur l'alphabet de taille finie  $(a \dots z, A \dots Z, \dots, 0 \dots 9, \dots, ;, \dots, !, \dots, ?)$ .

• Un texte en français est un signal sur l'alphabet  $A = \{ \text{mots français} \}$ . (85000 mots dans le Larousse)

Plutôt que de transmettre chaque symbole de l'alphabet sans codage est ce que on peut coder chaque symbole de manière à ce que les symboles fréquents aient des codes courts de manière à transmettre la même quantité d'informations en transmettant un message + coût ???

On se donne en modèle le message à transmettre c'est une probabilité sur  $A = (a_i)_{1 \leq i \leq k}$  et on se donne  $(p_{a_i})_{1 \leq i \leq k}$   $0 \leq p_{a_i} \leq 1$  et  $\sum_{i=1}^k p_{a_i} = 1$

On suppose qu'on essaye de transmettre

un message  $(X_1, X_2, \dots)$  où les  $X_i$  sont des v.a. indep à valeurs dans  $A$  et de loi donnée par  $p$  c'est à dire  $\mathbb{P}(X_j = a) = p_a$

Méthode naïve: si  $2^{p-1} < K \leq 2^p$

alors  $A$  en bijection avec une partie de  $\{0, 1\}^p$  on peut coder chaque lettre de  $A$  par un élément  $(a_i)$  de  $\{0, 1\}^p$  à un message  $x_1, \dots, x_n$  je transmet  $(x_1) \dots (x_n)$  le nombre moyen  $R_n$  de bits  $\{0, 1\}$  transmis vérifie  $R_n = np$ .

$$\text{Ici } R_n = n (\lfloor \log_2 K \rfloor + 1)$$

Et asymptotiquement le nombre moyen de bits transmis par lettre est

$$R = \lim_n \frac{R_n}{n} = \lfloor \log_2 K \rfloor + 1.$$

Objectif: Trouver d'autres codages pour lesquels  $R_n$  et  $R$  seront + petit.

Def: \*  $A = (a_i)_{1 \leq i \leq K}$  alphabet fini

\*  $p = (p_{a_i})_{1 \leq i \leq K}$  proba sur  $A$ .

\*  $p(x_1, \dots, x_n) = p_{x_1} \dots p_{x_n}$  proba d'un message  
 $= \mathbb{P}((X_1, \dots, X_n) = (x_1, \dots, x_n)).$

\* l'entropie de la distribution  $p$  :

$$H = - \sum_i p_{a_i} \log_2 p_{a_i}$$

Rem: Motivation pour cette définition :

Imaginons un ensemble  $\Omega = \{0, 1\}^N$   
muni de la proba uniforme  $\mathbb{P}(\{x\}) = \frac{1}{2^N}$ .

Je me fixe un élément  $x$  de  $\Omega$

Combien de bits d'information je dois transmettre pour vous dire de quel  $x$  il s'agit ?  $x = (x_1, x_2, \dots, x_N)$

Il faut vous transmettre  $N$  informations  
 $x_1 \in \{0, 1\}, x_2 \in \{0, 1\}, \dots, N = \log_2 |\Omega|$   
 $= -\log_2 \mathbb{P}(x)$

$A = \{x \in \Omega : x_1 = \varepsilon_1, \dots, x_p = \varepsilon_p\}$ .

Combien de bits d'information je dois transmettre pour préciser qu'un élément donné est dans  $A$  ?

Je dois dire  $x_1 = \varepsilon_1, \dots, x_p = \varepsilon_p$

donc  $p$  informations  $A = \{\varepsilon_1\} \times \{0, 1\} \times \dots \times \{\varepsilon_p\} \times \{0, 1\}^{N-p}$

$$\text{On } \mathbb{P}(A) = \frac{|A|}{|\Omega|} = \frac{2^{N-p}}{2^N} = 2^{-p}$$

Il faut  $-\log_2 \mathbb{P}(A)$  informations.

En général :  $-\log_2 \mathbb{P}(A)$  est le nombre de bits d'informations à transmettre pour dire qu'un élément est dans  $A$ .

Sur mon alphabet



$-\log_2 p_a$  = nbre d'information à transmettre pour dire qu'on a  $a$  (si  $a$  représentait une partie d'un  $\Omega = \{1, 2, \dots, N\}$ )

$$H = \sum_{a \in A} p_a (-\log_2 p_a) = \mathbb{E}[-\log_2 X]$$

= quantité d'information moyenne d'une lettre. où  $X$  de loi  $p$

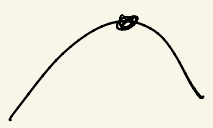
Rem:  $H = - \sum_{a \in A} p_a \log_2 p_a$

•  $H \geq 0$  SO car  $0 \leq p_a \leq 1$

•  $(x \log x)'' = (\log x + 1)' = \frac{1}{x} \geq 0$  si  $x > 0$

donc  $\psi(x) = -x \log x$  est concave et même strictement concave  $\frac{1}{x} > 0$  sur  $]0, 1[$ .

Donc:  $H = \sum_{a \in A} \psi(p_a) = K \sum_{a \in A} \frac{1}{K} \psi(p_a)$

  $\leq K \psi\left(\sum_{a \in A} \frac{1}{K} p_a\right)$  ( $K = |A|$ )

concavité  $= K \psi\left(\frac{1}{K} \sum_{a \in A} p_a\right)$

$= K \left(-\frac{1}{K} \log_2 \frac{1}{K}\right) = 1$

$= \log_2 K$

Donc on a  $0 \leq H \leq \log_2 K$ .

Le cas d'égalité a droite a lieu si  $p_a$  constant:  $p_a = \frac{1}{K} \forall a$ .

Donc borne de droite:

$$H = \log_2 K \Leftrightarrow \forall a \quad p_a = \frac{1}{K} \quad (\text{probas uniformes})$$

borne de gauche:

$$H = 0 \Leftrightarrow \forall a \quad p_a \log_2 p_a = 0$$

$$\Leftrightarrow \forall a \quad p_a \in \{0, 1\}$$

$$\Leftrightarrow \exists a_0 : p_a = \mathbb{1}_{a=a_0}$$

$$H = 0 \Leftrightarrow p \text{ est un dirac.}$$

$H$  mesure la quantité de désordre d'aléa du signal.

$H = 0 \Rightarrow$  signal déterministe

$H = \log_2 K \Rightarrow$  désordre maximal. loi uniforme.

Prop:  $X_1, X_2, \dots$  v.o. indep  $\in A$  de loi  $p$ .

$$P(X_1, \dots, X_n) = P_{X_1} \dots P_{X_n}$$

On:

$$(P_{x_1, \dots, x_n}) = \mathbb{1}(X_1 = x_1, \dots, X_n = x_n)$$

$$-\frac{1}{n} \log_2 P(X_1, \dots, X_n) \xrightarrow{n \rightarrow \infty} H \quad \text{avec proba } 1.$$

Rem: lorsque l'on regarde un long message

$X_1, \dots, X_n$  la proba de l'observer

est  $\approx 2^{-nH}$

Preuve: repose la loi forte des grands nombres

LFGN: Si  $Y_1, Y_2, \dots, Y_n$  sont intégrables <sup>et indep.</sup>  
 Alors  $\frac{1}{n} (Y_1 + \dots + Y_n) \xrightarrow{n \rightarrow \infty} E[Y_i]$  avec proba 1.

Ici:  $-\frac{1}{n} \log_2 P(X_1, \dots, X_n) \quad \triangle Pa \in \mathbb{R}$   
 $= \frac{1}{n} (Y_1 + \dots + Y_n) \quad P_{X_i} \text{ v.o.}$

avec  $Y_i = -\log_2 P_{X_i}$  les  $Y_i$  sont indep  
 et  $Y_i$  intégrable?  $E[|\log_2 P_{X_i}|]$   
 $\in \{\log_2 P_a \mid a \in A\}$

et  $\mathbb{P}(|\log_2 P_{X_i}| = \infty) = \sum_{P_a=0} \underbrace{\mathbb{P}(X_i=a)}_{P_a} = 0$

Donc  $E[|\log_2 P_{X_i}|]$  prend un nombre fini de valeurs fini avec proba 1 donc  $< \infty$

Donc  $-\frac{1}{n} \log_2 P(X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{\text{proba 1}} E[-\log_2 P_{X_i}] = -\sum P_a \log_2 P_a = H$

Novelté: Avec grande proba toutes les  $X_1, \dots, X_n$  que l'on observe avait proba  $\approx 2^{-nH}$  d'arriver.  
 On veut définir des suites "typiques":

$T_\varepsilon^n = \left\{ (x_1, \dots, x_n) \in A^n \mid H - \varepsilon \leq -\frac{1}{n} \log_2 P(x_1, \dots, x_n) \leq H + \varepsilon \right\}$   
 $= \left\{ (x_1, \dots, x_n) \in A^n \mid 2^{-n(H+\varepsilon)} \leq P(x_1, \dots, x_n) \leq 2^{-n(H-\varepsilon)} \right\}$

La proposition dit que l'on est "typique" avec

grande proba :

$$\mathbb{P}(X_1, \dots, X_n \in T_\varepsilon) \xrightarrow{n \rightarrow \infty} 1$$

(La proposition que'on a montrée est plus forte que ça).

En pratique : on ne voit pas n'importe quelle suite de  $A^n$  mais seulement celles de  $T_\varepsilon$

Théorème  $(1-\varepsilon)2^{n(H-\varepsilon)} \leq |T_\varepsilon^n| \leq 2^{n(H+\varepsilon)} \quad \forall \varepsilon$   
si  $n \geq N_\varepsilon$

Dem :  $\exists N_\varepsilon$  tel que  $n \geq N_\varepsilon$  :  
(cardinal)

$$1-\varepsilon \leq \mathbb{P}(X_1, \dots, X_n \in T_\varepsilon^n) \leq 1$$

Rappel :  $x_1, \dots, x_n \in T_\varepsilon^n \Leftrightarrow P_{x_1, \dots, x_n} \in [2^{-n(H+\varepsilon)}; 2^{-n(H-\varepsilon)}]$

$$\begin{aligned} 1 \geq \mathbb{P}(X_1, \dots, X_n \in T_\varepsilon^n) &= \sum_{x_1, \dots, x_n \in T_\varepsilon^n} \underbrace{\mathbb{P}(X_1, \dots, X_n = x_1, \dots, x_n)}_{P_{x_1, \dots, x_n}} \\ &\geq \sum_{x_1, \dots, x_n \in T_\varepsilon^n} 2^{-n(H+\varepsilon)} \\ &= |T_\varepsilon^n| 2^{-n(H+\varepsilon)} \end{aligned}$$

$$\text{Donc } |T_\varepsilon^n| \leq 2^{n(H+\varepsilon)} \quad \checkmark$$

$$\begin{aligned} 1-\varepsilon \leq \mathbb{P}(X_1, \dots, X_n \in T_\varepsilon^n) &\leq \sum_{x_1, \dots, x_n \in T_\varepsilon^n} P_{x_1, \dots, x_n} \\ &\leq |T_\varepsilon^n| 2^{-n(H-\varepsilon)}. \end{aligned}$$

$$\text{Donc } |T_\varepsilon^n| \geq (1-\varepsilon) 2^{n(H-\varepsilon)}.$$



Conséquence: codage: si on a un message  $x_1, \dots, x_n$  on le code de la manière suivante

\* si  $(x_1, \dots, x_n) \in T_\varepsilon^n$  on le code comme

$$0 \varphi(x_1, \dots, x_n)$$

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

$T_\varepsilon^n$  s'injecte dans  $\{0,1\}^{\lfloor n(H+\varepsilon) \rfloor + 1}$

$$\varphi: T_\varepsilon^n \rightarrow \{0,1\}^{\lfloor n(H+\varepsilon) \rfloor + 1}$$

\* si  $x_1, \dots, x_n \notin T_\varepsilon^n$

$$A^n \xrightarrow[\text{inj}]{\varphi} \{0,1\}^{\lfloor n \log_2 K \rfloor + 1}$$

(codage naïf)

on le code comme

$$1 \varphi(x_1, \dots, x_n)$$

Remarques: Facile à bécoder  $0 z_1 \dots z_m \rightarrow \varphi^{-1}(z_1 \dots z_m)$   
 $1 z_1 \dots z_m \rightarrow \varphi^{-1}(z_1 \dots z_m)$

•  $R_n$  = nbre moyen de bits produits.

$$= \underbrace{1}_{\text{over 1}} + (\lfloor n(H+\varepsilon) \rfloor + 1) \mathbb{P}(X_1, \dots, X_n \in T_\varepsilon^n)$$

$$+ n(\lfloor \log_2 K \rfloor + 1) \mathbb{P}(X_1, \dots, X_n \notin T_\varepsilon^n)$$

pour  $n$  grand

$$\leq 1 + (\lfloor n(H+\varepsilon) \rfloor + 1) \underbrace{1}_{\leq \varepsilon} + n(\lfloor \log_2 K \rfloor + 1) \varepsilon$$

$$\leq 2 + nH + nC\varepsilon$$

Donc le nombre moyen de bits par lettre est:

$$R = \lim_n \frac{R_n}{n} \leq H + C\varepsilon.$$

Le codage naïf bonnaît  $R \approx \log_2 K$ .

Ici on a  $R$  arbitrairement proche de  $H$  et on a vu  $H \leq \log_2 R$

On verra qu'on ne peut pas avoir de codage avec  $R < H$ .

→ codage optimal.

Mais \* codage pas explicite dépend de  $\mathcal{C}$  et \* il dépend de  $n$ . Si on rajoute des lettres au message il faut le coder depuis le début (pas un codage lettre à lettre)

Codage de Huffman :

• codage lettre à lettre : on construit

$$\mathcal{C} : A \longrightarrow \bigcup_{n \geq 0} \{0, 1\}^n$$

$\mathcal{C}(a)$  est le code de  $a$

$x_1 \dots x_n$  codé en  $\mathcal{C}(x_1) \dots \mathcal{C}(x_n)$

On ne peut pas forcément décoder un tel code : Si  $A = \{a, b, c\}$

$$\mathcal{C}(a) = 01 \quad \mathcal{C}(b) = 011 \quad \mathcal{C}(c) = 1$$

et que l'on reçoit 011 on ne sait pas si le message initial est ac ou b

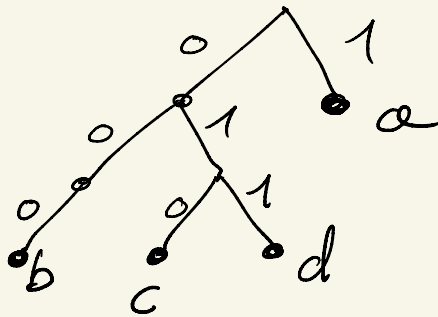
On demande à  $\mathcal{C}$  d'être "sans ambiguïté"

c'est à dire que  $\mathcal{C}(ai)$  n'est jamais

préfixe de  $\mathcal{C}(d_j)$  pour  $j \neq i$  ( $\mathcal{C}(a_j) \neq \mathcal{C}(a_i) \forall i, j$ )

Un tel codage peut être représenté sur un arbre binaire :

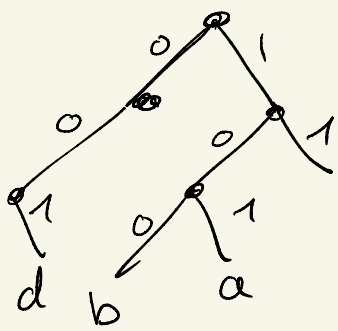
$$\begin{aligned} \text{Si } \mathcal{C}(a) &= 1 & \mathcal{C}(b) &= 000 \\ \mathcal{C}(c) &= 010 & \mathcal{C}(d) &= 011 \end{aligned}$$



On place chaque lettre dans l'arbre à la position de son code

La condition sans préfixe ( $\Rightarrow$ ) aucune lettre n'est descendante d'une autre

Réciproquement : si on se donne un arbre binaire avec des lettres sur les feuilles on en déduit un code sans préfixe :



on lit le code en descendant de la racine :

$$\begin{aligned} \mathcal{C}(a) &= 101 \\ \mathcal{C}(b) &= 100 \\ \mathcal{C}(c) &= 11 \\ \mathcal{C}(d) &= 001 \end{aligned}$$

Code de Huffman : on a un alphabet

$A = (a_i)_{1 \leq i \leq K}$  et une proba  $(P(a_i))_{1 \leq i \leq K}$

on classe les lettres de manière à avoir



$$P_{a_1} \leq P_{a_2} \leq \dots$$

on pose  $a_{12} = \{a_1, a_2\}$  on construit l'arbre de Huffman sur

$a_{12}, a_3, a_4, \dots$

avec  $P_{a_{12}} = P_{a_1} + P_{a_2}$

L'arbre final s'obtient en groupant  $a_1$  et  $a_2$



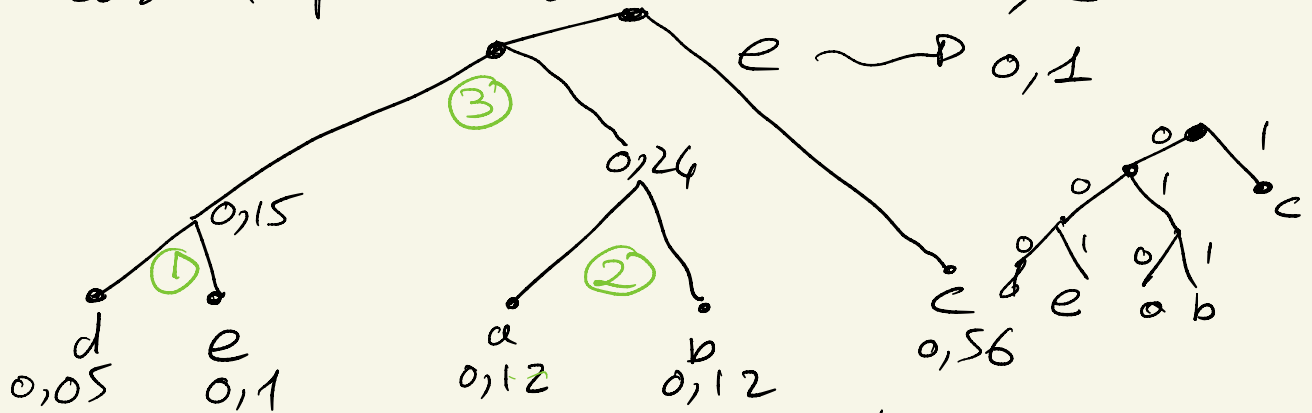
Puis on continue sur le nouvel alphabet

Ex:  $A = \{a; b; c; d; e\}$

$$P = \{0,12; 0,12; 0,56; 0,05; 0,1\}$$

Les deux + petits sont  $d \rightsquigarrow 0,05$

$e \rightsquigarrow 0,1$



$\varphi$ :

a	→	010
b		011
c		1
d		000
e		001

une lettre est beaucoup + fréquente "c"  
 $\rightsquigarrow$  elle est codée par un code très court.

Exo:  $A = \{M, P, R, U, Y, Z\}$

$$P_M = \frac{1}{16} = P_P = P_R = P_Z \quad P_U = \frac{1}{4} \quad P_Y = \frac{1}{2}$$

① Calculer H.



② Calculer le codage Huffman et l'arbre associé.

③ Calculer  $R_1 =$  nbre moyen de lettre du code d'une lettre aléatoire  
 $= E[\text{longueur du code de } X_1]$

④ Quel est le code de YUPT ?

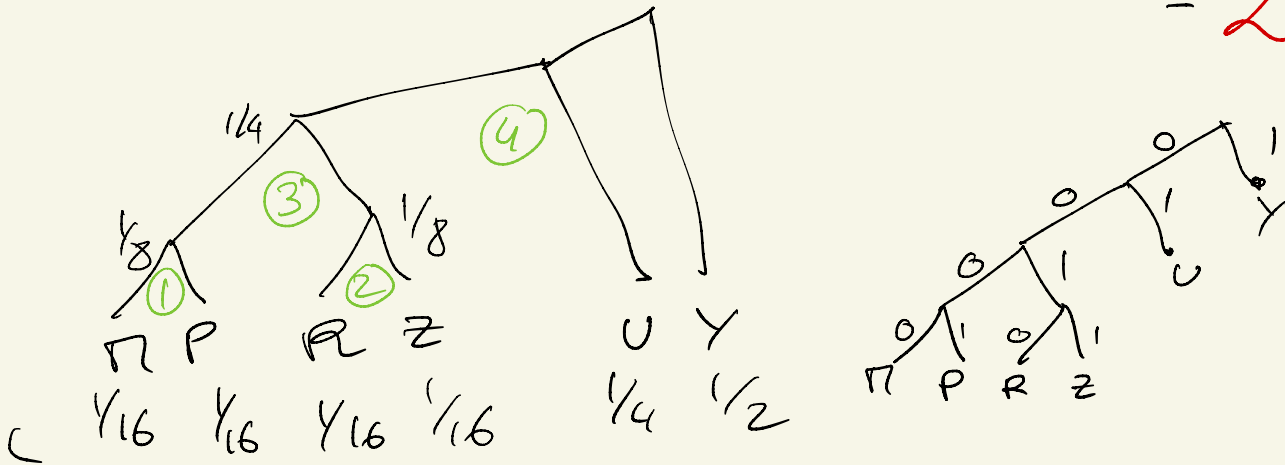
$$\log_2 \frac{1}{2^n} = -n$$

$$OH = -\sum p_a \log_2 p_a$$

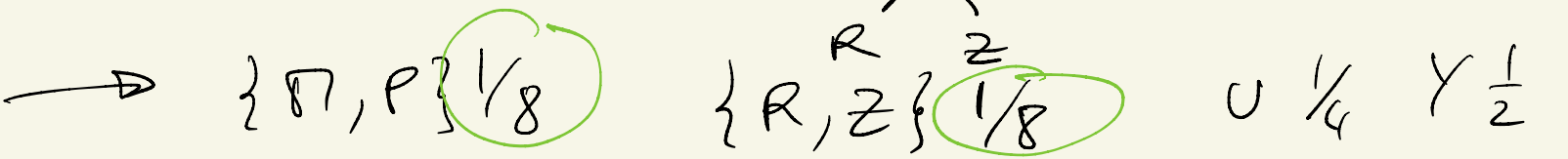
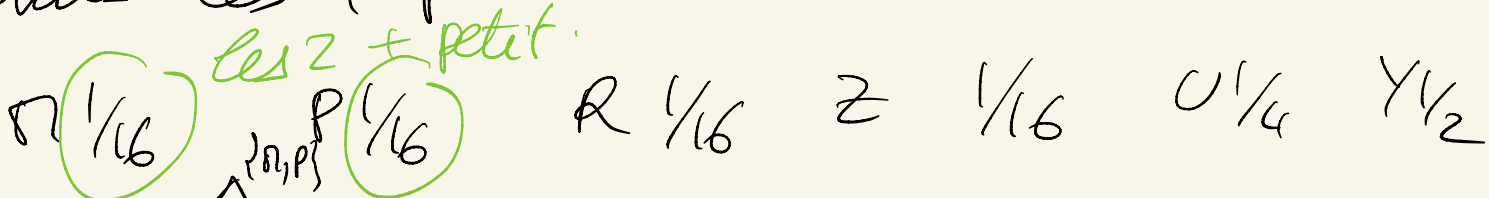
$$= -4 \frac{1}{16} \log_2 \frac{1}{2^4} - \frac{1}{4} \log_2 \frac{1}{2^2} - \frac{1}{2} \log_2 \frac{1}{2}$$

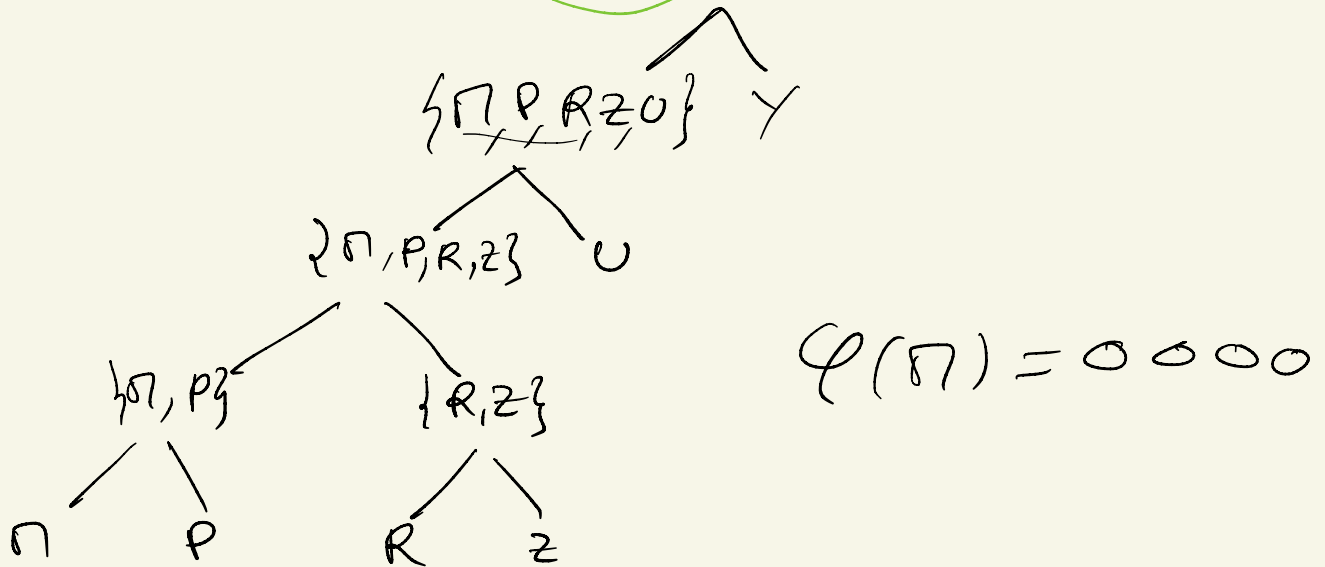
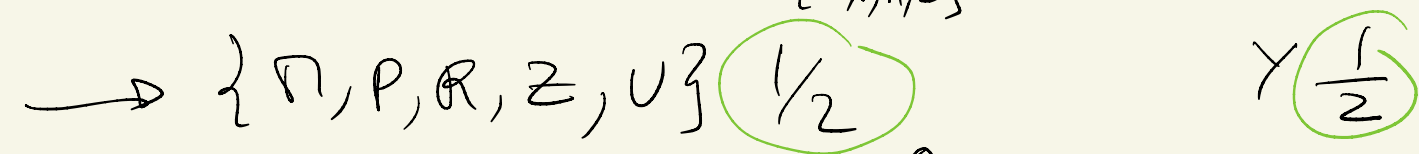
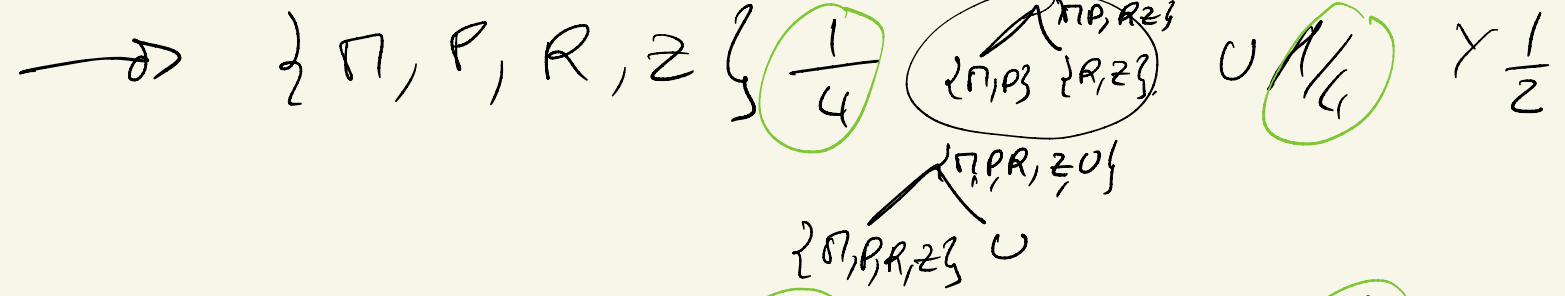
$$= -\frac{1}{4} (-4) - \frac{1}{4} (-2) - \frac{1}{2} (-1) = 1 + \frac{1}{2} + \frac{1}{2} = 2$$

②



A chaque étape je regroupe les 2 probas les + faibles





$\underline{\underline{R_1}} = E[\text{longueur des code}]$

$= 4 \cdot \mathbb{P}(X = \pi, P, R \text{ ou } Z)$   
 $+ 2 \cdot \mathbb{P}(X = U) + 1 \cdot \mathbb{P}(X = Y)$

$= 4 \times 4 \times \frac{1}{16} + 2 \times \frac{1}{4} + \frac{1}{2} = 2 = H$