

# Statistiques sur les lettres

Soit un document  $d$  :

- constitué de  $T$  symboles  $d[1], \dots, d[i], \dots$
- appartenant à l'alphabet  $A = \{\alpha_1, \dots, \alpha_K\}$  constitué de  $K$  symboles.

## Modèles probabilistes

Les modèles probabilistes interprètent les données de type texte comme étant générées par une distribution de probabilité  $P$  inconnue.

La distribution  $P$  définit le langage utilisé dans le texte. On ne s'intéresse pas au sens du message, on regarde seulement comment les symboles se répartissent dans les documents, leurs fréquences d'apparition, les régularités, ...

## Fréquence d'un symbole

Soit  $\alpha \in A$  un symbole de l'alphabet. On note  $P(X=\alpha)$  la fréquence d'apparition de ce symbole dans le langage  $\mathcal{L}$  considéré, soit :  $P(X=\alpha) = \frac{|\{\omega \in \Omega : X=\alpha\}|}{|\Omega|}$  où  $\Omega$  représente l'ensemble des productions de caractères.

On a par définition :  $\sum_{\alpha \in V} P(X=\alpha) = 1$

La fréquence empirique du symbole  $\alpha$  dans le document  $d$  est donnée par :



$$f_d(\alpha) = \frac{|\{i: d[i] = \alpha\}|}{|d|}$$

où  $|d|$  est le nombre de caractères dans le document.



## Fréquence des lettres en français

- Voir aussi : [Analyse Fréquentielle sur Wikipedia](#)

## Représentation vectorielle

On suppose que les caractères d'un langage  $\mathcal{L}$  donné sont numérotés de 1 à  $K$ , soit  $A_{\mathcal{L}} = \{\alpha_1, \dots, \alpha_k, \dots, \alpha_K\}$ .

On notera  $\mathbf{p}_{\mathcal{L}}$  le vecteur des fréquences des caractères dans un langage  $\mathcal{L}$  donné, où  $\mathbf{p}_{\mathcal{L}}(k)$  donne la fréquence du  $k^{\text{ème}}$  caractère.



**Exemple:**  $\mathbf{p}_{\text{Français}} = (0.0942, 0.0102, 0.0264, 0.0339,$

0.01587, 0.095, 0.0104, 0.0077, 0.0841, 0.0089, ...) où



- $p_1 = 0.0942$  est la fréquence de la lettre 'A',
- $p_2 = 0.0102$  est la fréquence d'apparition de la lettre 'B'
- etc.

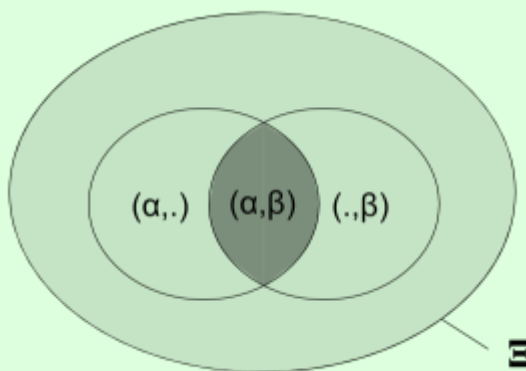
avec bien sûr :  $\sum_{k \in \{1, \dots, K\}} p_{\mathcal{L}}(k) = 1$

## Probabilité jointe

On s'intéresse maintenant aux fréquences d'apparition de couples de lettre successives.

Soient  $\alpha$  et  $\beta$  deux symboles de l'alphabet.

La probabilité jointe est définie comme :  $P(X=\alpha, Y=\beta) = \frac{|\{x \in X : (X,Y)=(\alpha,\beta)\}|}{|X|}$  où  $X$  est l'ensemble des productions de couples de caractères.



avec par définition:  $\sum_{(\alpha,\beta) \in A \times A} P(X=\alpha, Y=\beta) = 1$

La **probabilité jointe empirique** est donnée par :



$$f_d(\alpha, \beta) = \frac{|\{i: d[i] = \alpha, d[i+1] = \beta\}|}{|d|-1}$$

- Les séquences de deux caractères sont classiquement appelées des *bigrammes*.
- On définit de même les *trigrammes* comme les séquences de trois caractères
- etc.

## Représentation matricielle

On notera  $P_{\mathcal{L}}$  la matrice des fréquences des bigrammes dans un langage  $\mathcal{L}$  donné, où  $P_{ij}$  donne la fréquence du bigramme  $(\alpha_i, \alpha_j)$ .

**Exemple:**  $P_{\text{Français}} = 10^{-5} \times \left( \begin{array}{cccc} 1.5 & 116.8 & 199.1 & \dots \\ 62.8 & 1.6 & 0.14 & \dots \\ 184.8 & 0 & 52.4 & \dots \end{array} \right)$



où

- $P_{11} = 1.5 \times 10^{-5}$  est la fréquence du bigramme 'AA',
- $P_{12} = 116.8 \times 10^{-5}$  est la fréquence d'apparition du bigramme 'AB'
- etc.

avec bien sûr :  $\sum_{(i,j) \in \{1,\dots,K\}^2} P_{ij} = 1$



voir [comptage des bigrammes en français](#)

## Corpus de documents

Soit  $B$  un corpus de documents, constitué de  $n$  documents.



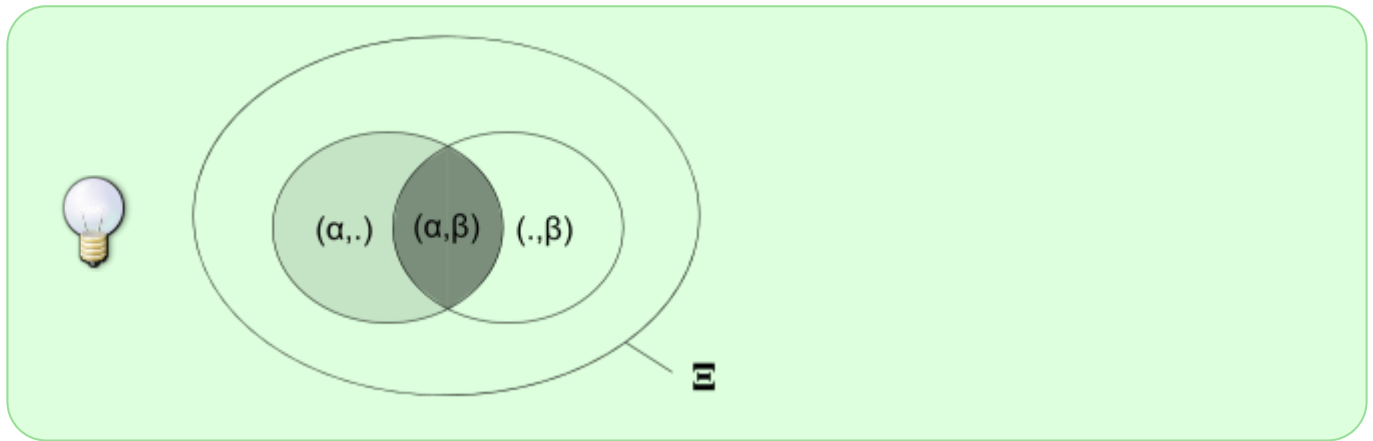
La fréquence empirique du symbole  $\alpha$  dans le corpus  $B$  est donnée par~:  $f_B(\alpha) = \frac{|\{(i,j): d_i \in B, d_{i[j]} = \alpha\}|}{|B|}$  où  $|B|$  est le nombre total de caractères dans le corpus.

La fréquence jointe du couple  $(\alpha, \beta)$  est donnée par  $f_B(\alpha, \beta) = \frac{|\{(i,j): d_i \in B, (d_{i[j]}, d_{i[j+1]}) = (\alpha, \beta)\}|}{|B|-n}$

## Probabilité conditionnelle

La **probabilité conditionnelle** du caractère  $\beta$  étant donné le caractère précédent  $\alpha$  est définie comme :

$$P(Y = \beta \mid X = \alpha) = \frac{|\{x_i \in X_i : (X, Y) = (\alpha, \beta)\}|}{|\{x_i \in X_i : X = \alpha\}|}$$



qui se calcule empiriquement comme :

$$f_d(\beta|\alpha) = \frac{|\{i:d[i] = \alpha, d[i+1] = \beta\}|}{|\{j:d[j] = \alpha\}|}$$

- La probabilité  $P(.|\alpha_i)$  se représente sous forme vectorielle~:  $\boldsymbol{\mu}_i = (P(\alpha_1|\alpha_i), P(\alpha_2|\alpha_i), \dots)$  où  $\alpha_1$  est le premier caractère de l'alphabet,  $\alpha_2$  le deuxième etc, avec  $\sum_j \boldsymbol{\mu}_{ij} = 1$



- L'ensemble des probabilités conditionnelles  $P(.|.)$  peut se représenter sous une forme matricielle~:

$$M = \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \\ \dots \end{pmatrix} = \begin{pmatrix} P(\alpha_1|\alpha_1) & P(\alpha_2|\alpha_1) & P(\alpha_3|\alpha_1) & \dots \\ P(\alpha_1|\alpha_2) & P(\alpha_2|\alpha_2) & P(\alpha_3|\alpha_2) & \dots \\ \dots & \dots & \dots & \dots \end{pmatrix}$$

Sachant que  $P(\alpha) = \sum_{\beta \in A} P(\alpha, \beta)$ , on a :  $\boldsymbol{\mu}_i = \frac{P_{i,:}}{p_i}$

Soit en français :

$$M_{\text{Français}} = \begin{pmatrix} 0.0016 & 0.0124 & 0.0211 & \dots \\ 0.0615 & 0.0016 & 0.0001 & \dots \\ 0.0700 & 0.0000 & 0.0198 & \dots \end{pmatrix}$$

où :

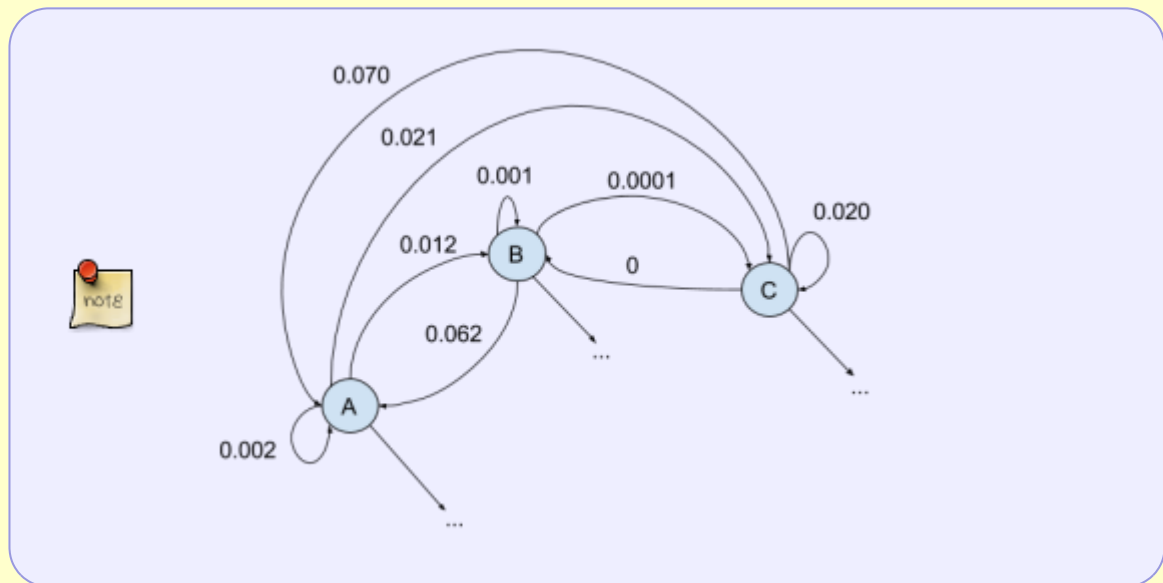


- $M_{11}$  est la probabilité de voir un 'A' suivre un 'A'
- $M_{12}$  est la probabilité de voir un 'B' suivre un 'A'
- etc.



La matrice des probabilités conditionnelles  $M$  permet de définir un **modèle génératif** de langage inspiré des **processus aléatoires de Markov**:

- La production d'un mot ou d'un texte est modélisée comme un parcours aléatoire sur une chaîne de Markov définie par la matrice de transitions  $M$ .
- La fréquence d'apparition des lettres est modélisée comme la mesure stationnaire de la chaîne de Markov, autrement dit le vecteur de probabilité vérifiant :  $\pi M = \pi$



## Comparer les langues

On considère deux langues  $\mathcal{L}_1$  et  $\mathcal{L}_2$  utilisant le même alphabet. La différence de fréquence des caractères dans ces deux langages permet de les distinguer. Il est ainsi possible de définir une distance entre deux langages basée sur la distance Euclidienne entre les vecteurs de fréquence empirique des caractères dans les deux langages.

Une autre approche consiste à utiliser la *théorie de l'information* pour comparer deux langues.

L'information apportée par la lecture du symbole  $\alpha$  est définie comme :  $I(\alpha) = -\log_2(P(X = \alpha))$  où  $p_\alpha = P(X = \alpha)$  est la fréquence d'apparition de ce symbole dans la langue considérée.

Si le symbole  $\alpha$  est "rare" ( $p_\alpha$  petit), l'information qu'il apporte est élevée. Si le symbole est fréquent, l'information qu'il apporte est faible.

L'entropie d'un langage est définie comme l'espérance de l'information apportée par un caractère.  $H(\mathcal{L}) = E_X(I(\alpha)) = -E_X(\log_2(P(X = \alpha)))$  i.e.  $H(\mathcal{L}) = -\sum_{k \in \{1, \dots, K\}} P(X = \alpha_k) \log_2(P(X = \alpha_k))$

L'entropie représente l' "imprévisibilité" d'une production de symboles. Une entropie faible indique que la séquence est très prévisible, une entropie élevée indique une séquence très imprévisible.

Pour comparer deux langues  $\mathcal{L}_1$  et  $\mathcal{L}_2$ , on utilise dans ce cadre la *divergence de Kullback-Leibler* définie comme:

$$D(\mathcal{L}_1 || \mathcal{L}_2) = \sum_{k \in \{1, \dots, K\}} P_1(X = \alpha_k) \log_2 \left( \frac{P_1(X = \alpha_k)}{P_2(X = \alpha_k)} \right)$$

où  $P_1$  désigne la distribution des symboles du langage  $\mathcal{L}_1$  et  $P_2$  la distribution des symboles du langage  $\mathcal{L}_2$ .

La divergence de K-L représente l'espérance, dans le langage  $\mathcal{L}_1$ , de la différence d'information apportée par un même symbole. Elle vaut 0 si les deux distributions sont identiques.

From:

<https://wiki.centrale-med.fr/informatique/> - **WiKi informatique**

Permanent link:

[https://wiki.centrale-med.fr/informatique/public:algo-txt:statistiques\\_sur\\_les\\_lettres](https://wiki.centrale-med.fr/informatique/public:algo-txt:statistiques_sur_les_lettres)
Last update: **2016/03/14 09:43**