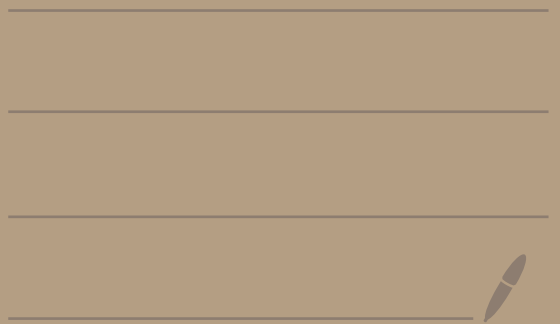


U E Modélisation :

Topic 2 :



Prédiction de mots par chaînes de Markov.

ou « comment écrire comme votre
auteur préféré »

Dans cette leçon qui peut être vue comme une introduction au NLP (Natural Language Processing), nous allons faire de la prédiction de mots.

Idee: Voir un texte comme une réalisation de chaîne de Markov sur $E = \{ \text{Mots} \}$

En fait Andreï Andreïevitch Markov (1856 - 1922) a conçu les chaînes de Markov avec l'exemple du TEXTE - juste avec deux états $E = \{ V = \text{"voyelle"}, C = \text{"consonne"} \}$.

I Rappel: Pour ce cours, nous avons uniquement besoin de 2 notions de la séance précédente:

(1) Définition des chaînes de Markov. (CM)

Si $E = \{ 1, 2, \dots, n \}$ espace d'états fini.

$P = (P_{ij})_{1 \leq i, j \leq n}$ matrice stochastique ($P_{ij} \geq 0$)

Alors la CM avec $\left\{ \begin{array}{l} \text{espace d'états } E \\ \text{matrice de transit } P \end{array} \right.$ $\sum_{j=1}^n P_{ij} = 1 \quad \forall i$

est le processus $X = (X_n; n \in \mathbb{N})$ donné par

$\left\{ \begin{array}{l} X_0 \in E \text{ (éventuellement aléatoire)} \\ P_{x,y} = \mathbb{P}(X_{n+1}=y | X_n=x) = \mathbb{P}(X_{n+1}=y | X_n=x, X_{n-1}, \dots, X_0) \end{array} \right.$

Moralement c'est l'exemple le plus simple après une suite iid. On parle souvent de dépendance faible car la dépendance est uniquement en l'étape précédente.

(\neq AR(p) ou MA(q) qui seront vues en séries chronologiques)

(2) Théorème ergodique de Birkhoff:

ou "loi forte des grands nombres"

Si X irréd, ~~apériodique~~

$$\text{Alors } \frac{1}{T} \sum_{t=0}^{T-1} f(X_t) \xrightarrow[T \rightarrow \infty]{P.S.} \mathbb{E}_{\pi}(f(X)) = \sum_{x \in E} \pi_x f(x)$$

où π unique mesure invariante ($\pi P = \pi$)

Problème: $E = \{ \text{Mots dans un texte} \}$

$P =$ Matrice de transitions pour la chaîne de Markov dont UNE réalisation est notre texte.

Comment estimer ce P ?

↳ **Statistique**

II L'exemple historique de Markov (Andrei)

Andrei Markov était en désaccord avec Pavel Nekrasov sur le fait que l'indépendance est nécessaire pour la loi (faible) des grands nombres.

A cet effet, dans un papier de 1913, Markov a choisi une phrase de longueur $T = 20\,000$ lettres du livre "Eugene Onegin" de Pushkin.

Il a estimé une chaîne de Markov sur $E = \{V, C\}$

avec

$$\hat{P}_T = \begin{matrix} & \begin{matrix} V & C \end{matrix} \\ \begin{matrix} V \\ C \end{matrix} & \begin{pmatrix} 0,128 & 0,872 \\ 0,663 & 0,337 \end{pmatrix} \end{matrix}$$

Puis il a calculé $\pi = (0,432; 0,568)$ où $\pi P = \pi$,
 et prouvé la loi (faible) des grands nombres
 / le thm de CV à l'équilibre (version faible de Birkhoff).

Question: Comment a-t-il trouvé ce \hat{P}_T ?

↳ Il a pris ce texte de longueur $T = 20\,000$.

↳ Il a soigneusement noté la réalisat^o correspondante de la CM X de matrice P :

Ex: Eugene Onegin.

VVCVGV VCVCVC

Converge (en Proba ou p.s)

↳ Calcul: $\hat{P}_T = \begin{pmatrix} \hat{\alpha}_T & 1 - \hat{\alpha}_T \\ 1 - \hat{\beta}_T & \hat{\beta}_T \end{pmatrix} \xrightarrow{T \rightarrow +\infty} P = \begin{pmatrix} \alpha & 1 - \alpha \\ 1 - \beta & \beta \end{pmatrix}$

où $\hat{\alpha}_T = \frac{\#\{0 \leq t \leq T-1 \mid X_t = X_{t+1} = V\}}{\#\{0 \leq t \leq T-1 \mid X_t = V\}} \xrightarrow{?} \alpha$

$\hat{\beta}_T = \frac{\#\{0 \leq t \leq T-1 \mid X_t = X_{t+1} = C\}}{\#\{0 \leq t \leq T-1 \mid X_t = C\}} \xrightarrow{?} \beta$

↳ Pourquoi cela fonctionne-t-il ?

Théorème: Soit $(X_t, t \in \mathbb{N})$ une chaîne de Markov sur E .

Alors $[\hat{P}_T]_{ij} := \frac{m_{ij}}{\sum_{k=1}^n m_{ik}}$ où $m_{ij} = \#\{0 \leq t \leq T-1 \mid (X_t, X_{t+1}) = (i, j)\}$

donne • une matrice stochastique \hat{P}_T .

• \hat{P}_T est un estimateur consistant de P

$$\left(\hat{P}_T \xrightarrow[T \rightarrow +\infty]{\mathbb{P}/\text{ps}} P \right)$$

$$\left(\sum_k m_{ik} = \#\{0 \leq t \leq T-1 \mid X_t = i\} \right)$$



La preuve probabiliste par le théorème ergodique de Birkhoff.

Exercice guidé dans le cas Markov

$$P = \begin{pmatrix} \alpha & 1-\alpha \\ 1-\beta & \beta \end{pmatrix} \text{ avec } \begin{cases} \alpha = \mathbb{P}(X_{t+1} = V \mid X_t = V) \\ \beta = \mathbb{P}(X_{t+1} = C \mid X_t = C) \end{cases}$$

① Par Birkhoff :

$$\frac{1}{T} \#\{0 \leq t \leq T-1 \mid X_t = V\} \xrightarrow[T \rightarrow +\infty]{} ?$$

$$\frac{1}{T} \#\{0 \leq t \leq T-1 \mid X_t = C\} \xrightarrow[T \rightarrow +\infty]{} ?$$

$$\mathbb{E}_{\pi} \mathbb{1}_V(X) = \frac{\mathbb{P}(X=V)}{\pi} = \pi_V$$

$$\pi_C$$

② Que reste-t-il à montrer pour que

$$\hat{\alpha}_T = \frac{\frac{1}{T} \#\{0 \leq t \leq T-1 \mid X_t = X_{t+1} = V\}}{\frac{1}{T} \#\{0 \leq t \leq T-1 \mid X_t = V\}} \xrightarrow[T \rightarrow +\infty]{} \alpha = P_{VV}$$

il reste

$$\frac{1}{T} \#\{0 \leq t \leq T-1 \mid X_t = X_{t+1} = V\} \xrightarrow[T \rightarrow +\infty]{} \pi_V P_{VV}$$

$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{1}_{\{X_t, X_{t+1}\}}$ // A besoin d'instants successifs! Pas Birkhoff

Invitation à se demander si (X_t, X_{t+1}) chaîne de Markov.

Sondage: OUI ou NON ?
 (1) (2)

Lemme: $(X_t, X_{t+1}) ; t \in \mathbb{N}$ est une chaîne de Markov de matrice de transition $Q = (Q_{(i,j), (k,e)})_{\substack{1 \leq i,j \leq n \\ 1 \leq k,e \leq n}}$

$$Q_{(i,j), (k,e)} = \mathbb{P}((X_{t+1}, X_{t+2}) = (k,e) \mid (X_t, X_{t+1}) = (i,j), \dots)$$

$$= \mathbb{1}_{\{j=k\}} \mathbb{P}_{R,e}$$

et de mesure invariante π_Q où

$$\pi_Q(k,e) = \pi_{\mathbb{P}}(k) \mathbb{P}_{R,e}$$

Preuve:

$$Q_{(i,j), (k,e)} = \mathbb{P}((X_{t+1}, X_{t+2}) = (k,e) \mid (X_t, X_{t+1}) = (i,j), \dots)$$

$$= \mathbb{1}_{\{j=k\}} \mathbb{P}(X_{t+1} = k, X_{t+2} = e \mid X_{t+1} = k, X_t = i, \dots)$$

Markov

$$= \mathbb{1}_{\{j=k\}} \mathbb{P}(X_{t+2} = e \mid X_{t+1} = k)$$

$$= \mathbb{1}_{\{j=k\}} \mathbb{P}_{R,e}$$

$$\pi_Q(k,e) \stackrel{\text{def}}{=} \pi_{\mathbb{P}}(k) \mathbb{P}_{R,e} \stackrel{?}{=} (\pi_Q Q)(k,e)$$

$$= \sum_{i,j} \pi_Q(i,j) Q(i,j), (k,e)$$

$$\begin{aligned}
\text{juste avant } \overline{\mathbb{P}} &= \sum_{i,j} \pi_Q(i,j) \mathbb{1}_{\{j=R\}} \mathbb{P}_{R,e} \\
&= \sum_i \pi_Q(i,R) \mathbb{P}_{R,e} = \left[\sum_i \pi_Q(i,R) \right] \mathbb{P}_{R,e} \\
&\stackrel{\text{def } \pi_Q}{=} \underbrace{\left(\sum_i \pi_P(i) \mathbb{P}_{i,R} \right)}_{\pi_P(R)} \mathbb{P}_{R,e} \\
&\qquad\qquad\qquad \text{def de } \pi_P \\
&\qquad\qquad\qquad \text{mesure invariante} \\
&= \pi_P(R) \mathbb{P}_{R,e} \\
&= \pi_Q(R,e) \quad \square
\end{aligned}$$

Ainsi en appliquant Birkhoff à (X_t, X_{t+1}) (~~*)~~)

$$\begin{aligned}
&\frac{1}{T} \# \{0 \leq t \leq T-1 \mid X_t = X_{t+1} = v\} \\
&= \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{1}_{\{v,v\}}(X_t, X_{t+1}) \\
&\xrightarrow[T \rightarrow +\infty]{(*)} \mathbb{E}_{\pi_Q} \mathbb{1}_{\{v,v\}}(X,Y) = \mathbb{P}_{\pi_Q}((X,Y) = (v,v)) \\
&= \pi_Q(v,v) \\
&\stackrel{\text{def}}{=} \pi_v \mathbb{P}_{vv} \quad \text{CQFD}
\end{aligned}$$

Sondage: Qui avait repéré ou était utilisée

la formule $[\hat{P}_T]_{ij} = \frac{n_{ij}}{\sum_{k=1}^m n_{ik}}$ dans le TP ?

- ① J'avais repéré.
- ② Non je n'avais repéré.
- ③ Je ne comprends pas la question.

Au programme aujourd'hui: Faire les preuves du thm
i.e. $\hat{P}_T \xrightarrow[T \rightarrow +\infty]{PS/P} P$ sur la vraie matrice de transition.

Rappel: En statistique, on dit que \hat{P}_T est un estimateur consistant.

III la preuve probabiliste par le thm. ergodique de Birkhoff
(suite et fin).

$$[\hat{P}_T]_{ij} = \frac{n_{ij}}{\sum_{k=1}^m n_{ik}} = \frac{\frac{1}{T} \#\{0 \leq t \leq T-1 \mid (X_t, X_{t+1}) = (i, j)\}}{\frac{1}{T} \#\{0 \leq t \leq T-1 \mid X_t = i\}}$$

? $\xrightarrow[T \rightarrow +\infty]{} P_{ij}$

A montrer pour (i, j) fixé

Comme pour le cas 2×2 avec les consonnes - voyelles de "Eugene Onegin", le dénominateur est traité par le thm de Birkhoff (en supposant irréductible + aperiodique).

$$\frac{1}{T} \# \{0 \leq t \leq T-1 \mid X_t = i\} \xrightarrow[T \rightarrow +\infty]{\text{Birkhoff (p.s)}} \frac{E_{\pi}[\mathbb{1}_{\{i\}}(X)]}{\pi(i)}$$

où π mesure invariante

Pour le numérateur, on écrit:

$$\frac{1}{T} \# \{0 \leq t \leq T-1 \mid (X_t, X_{t+1}) = (i, j)\} \xrightarrow[T \rightarrow +\infty]{\text{p.s.}} \pi_Q(i, j)$$

par Birkhoff appliqué à $(X_t, X_{t+1}) ; t \in \mathbb{N}$ est une chaîne de Markov

où π_Q probabilité invariante de $Y = ((X_t, X_{t+1}) ; t \in \mathbb{N})$

Rappel (Cours 2): $\pi_Q(\mathbb{R}, e) = \pi_{\mathbb{R}}(\mathbb{R}) \mathbb{P}_{\mathbb{R}, e}$

Donc $\frac{1}{T} \# \{0 \leq t \leq T-1 \mid (X_t, X_{t+1}) = (i, j)\} \xrightarrow[T \rightarrow +\infty]{\text{p.s.}} \pi(i) \mathbb{P}_{ij}$

Pour quotient:

$$[\hat{P}_T]_{ij} = \frac{\frac{1}{T} \# \{0 \leq t \leq T-1 \mid (X_t, X_{t+1}) = (i, j)\}}{\frac{1}{T} \# \{0 \leq t \leq T-1 \mid X_t = i\}} \xrightarrow[T \rightarrow +\infty]{\text{p.s.}} \frac{\pi(i) \mathbb{P}_{ij}}{\pi(i)} = \mathbb{P}_{ij}$$

CQFD.

Exo: Irréd. $\Rightarrow \forall i \in E, \pi(i) > 0$. □

Commentaire final: Sondage. Maintenant que vous avez vu la preuve et les arguments qu'elle utilise (CM, Birkhoff, Birkhoff pour les paires, Ergodicité),

me pensez-vous pas que la formule était plutôt intuitive et n'avait pas besoin de tout cela ?

(1) Oui **cette** preuve est un oraculé. (16)

(2) Non la preuve est nécessaire et la formule n'est pas intuitive. (12)

Commentaire final (bis).

$$\text{AR}(1): X_{t+1} = a X_t + \varepsilon_t \quad X = (X_t; t \in \mathbb{Z})$$

X a la propriété de Markov:

$$\mathbb{P}(X_{t+1} | \underbrace{X_t, X_{t-1}, \dots}_x) = \mathbb{P}(ax + \varepsilon).$$

$$\text{AR}(2): X_{t+1} = a_0 X_t + a_1 X_{t-1} + \varepsilon_t$$

X n'est pas Markov.

Mais $Y = ((X_t, X_{t+1}); t \in \mathbb{Z})$ est Markov.

$$\text{En effet: } Y_t = \begin{bmatrix} X_t \\ X_{t+1} \end{bmatrix} \quad Y_t$$

$$\begin{aligned} Y_{t+1} &= \begin{bmatrix} X_{t+1} \\ X_{t+2} \end{bmatrix} = \begin{bmatrix} X_{t+1} \\ a_0 X_{t+1} + a_1 X_t + \varepsilon_t \end{bmatrix} = \begin{bmatrix} 0 \\ \varepsilon_t \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ a_1 & a_0 \end{bmatrix} \begin{bmatrix} X_t \\ X_{t+1} \end{bmatrix} \\ &= \begin{bmatrix} 0 \\ \varepsilon_t \end{bmatrix} + \begin{bmatrix} 0 & 1 \\ a_1 & a_0 \end{bmatrix} Y_t \end{aligned}$$

Plus généralement, dans le $\text{AR}(p)$, une "fenêtre" de longueur p peut avoir la propriété de Markov.

$$Y^{(p)} = ((X_t, X_{t+1}, \dots, X_{t+p}); t \in \mathbb{Z}) \text{ Markov}$$

→ Moralement, notre TD2 était une histoire de processus autoregressif (AR) sur l'espace des mots.

IV La preuve du statisticien (avec une bonne dose d'optimisation)

« Méta-théorie de la statistique » :

Il faut considérer l'estimateur de maximum de vraisemblance (MLE: Maximum Likelihood Estimator).

Exercice pour se rafraîchir la mémoire.

(X_1, X_2, \dots, X_n) un échantillon iid de temps d'attente à l'arrêt d'un bus.

On suppose une loi exponentielle de densité

$$f_{\theta}(x) = \theta e^{-\theta x} \mathbb{1}_{x > 0} \quad \theta \text{ paramètre inconnu.}$$

Rappel: $E X = \frac{1}{\theta} \Rightarrow \hat{\theta}_n = \frac{1}{\hat{m}_n} = \frac{n}{\sum_{i=1}^n X_i}$

Calculer le MLE: $\hat{\theta}_n^*$ estimateur consistant.

$$\begin{aligned} \hat{\theta}_n^* &= \underset{\theta}{\operatorname{Argmax}} \ln L(\theta) && \text{(loi des grands nombres).} \\ &\stackrel{\text{L}}{\Rightarrow} \text{vraisemblance de} && \text{l'échantillon } (X_1, \dots, X_n) \\ &= \underset{\theta}{\operatorname{Argmax}} \log L_n(\theta) && = (x_1, \dots, x_n) \end{aligned}$$

"vraisemblance" = Probabilité / densité de voir ce que je vois.

$$\Rightarrow \ln L(\theta) = \text{Densité de } (X_1, X_2, \dots, X_n)$$

independance \downarrow

$$= f_{\theta}(x_1) f_{\theta}(x_2) \dots f_{\theta}(x_n)$$

$$\begin{aligned}
 \leadsto \log L_m(\theta) &= \sum_{i=1}^m \log f_{\theta}(x_i) \\
 &= \sum_{i=1}^m \log \left(\theta e^{-\theta x_i} \mathbb{1}_{\{x_i > 0\}} \right) \\
 &= \sum_{i=1}^m \left[\log \theta + \underbrace{\log e^{-\theta x_i}}_{-\theta x_i} + \underbrace{\log \mathbb{1}_{\{x_i > 0\}}}_{-\infty} \right] \\
 &= m \log \theta - \theta \sum_{i=1}^m x_i - \infty \mathbb{1}_{\{x_1 \leq 0 \vee x_2 \leq 0 \vee \dots \vee x_m \leq 0\}}
 \end{aligned}$$

Donc $\hat{\theta}_m^* = \underset{\theta}{\operatorname{Argmax}} \underbrace{m \log \theta - \theta \sum_{i=1}^m x_i}_{\varphi(\theta)}$

$$\varphi'(\theta) = \frac{m}{\theta} - \sum_{i=1}^m x_i \stackrel{\varphi(\theta)}{=} 0$$

$$\Leftrightarrow \frac{m}{\theta} = \sum_{i=1}^m x_i \quad \Leftrightarrow \theta = \frac{m}{\sum_{i=1}^m x_i}$$

Donc $\hat{\theta}_m^* = \text{MLE} = \frac{m}{\sum_{i=1}^m x_i} = \hat{\theta}_m$ ■

Revenons à notre chaîne de Markov.

On observe en échantillon $(x_0, x_1, x_2, \dots, x_T)$

↳ Paramètre: $\theta = \mathbb{I} = (x_0, x_1, x_2, \dots, x_T)$

↳ Calcul de la vraisemblance

$$L = L(\mathbb{I}) = \mathbb{P}(X_T = x_T, X_{T-1} = x_{T-1}, \dots, X_0 = x_0)$$

$$= \mathbb{P}(X_T = x_T \mid X_{T-1} = x_{T-1}, \dots, X_0 = x_0) \stackrel{\text{Markov}}{\downarrow}$$

$$\times \mathbb{P}(X_{T-1} = x_{T-1}, \dots, X_0 = x_0) \quad \mathbb{P}(X_T = x_T \mid X_{T-1} = x_{T-1})$$

$$= \mathbb{P}_{x_{T-1}, x_T}$$

$$= \mathbb{P}_{x_{T-1}, x_T}$$

$$\times \mathbb{P}(X_{T-1} = x_{T-1}, \dots, X_0 = x_0)$$

= ... Réurrence immédiate

$$= \mathbb{P}_{x_{T-1}, x_T} \times \mathbb{P}_{x_{T-2}, x_{T-1}} \times \dots \times \mathbb{P}_{x_1, x_0}$$

$$= \left(\prod_{i=0}^{T-1} \mathbb{P}_{x_i, x_{i+1}} \right) \times \mathbb{P}(X_0 = x_0) \quad \left| \quad \mathbb{P}(X_0 = x_0) \right.$$

↳ Ré-écriture :

$$L(\mathbb{P}) = \mathbb{P}(X_0 = x_0) \times \prod_{1 \leq i, j \leq m} \mathbb{P}_{ij}^{\#\{(i,j) \text{ apparait dans } (x_0, x_1, \dots, x_T)\}}$$

$$= \mathbb{P}(X_0 = x_0) \times \prod_{1 \leq i, j \leq m} \mathbb{P}_{ij}^{n_{ij}}$$

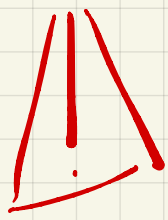
Donc MLE est donné par

$$\hat{\mathbb{P}}_T = \operatorname{Argmax}_{\mathbb{P}} \log L(\mathbb{P})$$

$$= \operatorname{Argmax}_{\mathbb{P}} \log \mathbb{P}(X_0 = x_0) + \sum_{1 \leq i, j \leq m} \log \mathbb{P}_{ij}^{n_{ij}}$$

$$\Rightarrow \hat{\mathbb{P}}_T = \operatorname{Argmax}_{\mathbb{P}} \underbrace{\sum_{1 \leq i, j \leq m} n_{ij} \log \mathbb{P}_{ij}}_{L(\mathbb{P})}$$

$$L(\mathbb{P})$$



Il faut optimiser sur un
espace de matrices \mathcal{P}

de taille $m \times n$

sous contrainte $\sum_{j=1}^m P_{ij} = 1 \quad \forall i$

↑
matrice
stochastique

↳ Calcul par la méthode
des multiplicateurs de Lagrange :

Poser $\mathcal{L}_c(\mathcal{P}, \lambda) = \mathcal{L}(\mathcal{P}) - \sum_{i=1}^m \lambda_i \left(\sum_{j=1}^m P_{ij} - 1 \right)$

↑
loss précédente

puis optimiser comme
si \mathcal{P}
et $\lambda = (\lambda_1, \dots, \lambda_m)$
sont libres / sans contraintes.

multiplicateur / contrainte
de Lagrange pour
la contrainte i .

(\mathcal{P}, λ) point critique

$$\Leftrightarrow \begin{cases} \forall i, \frac{\partial \mathcal{L}_c}{\partial \lambda_i} = 0 \\ \forall ij, \frac{\partial \mathcal{L}_c}{\partial P_{ij}} = 0 \end{cases} \Leftrightarrow \begin{cases} \sum_{j=1}^m P_{ij} - 1 = 0 \quad \forall i \\ \forall ij, \frac{m_{ij}}{P_{ij}^2} - \lambda_i = 0 \end{cases}$$

$$\Leftrightarrow \begin{cases} \forall ij, P_{ij} = \frac{m_{ij}}{\lambda_i} \\ \forall i, 1 = \sum_{j=1}^m P_{ij} = \left(\sum_{j=1}^m m_{ij} \right) / \lambda_i \Rightarrow \lambda_i = \sum_{j=1}^m m_{ij} \end{cases}$$

Donc

$$F_{ij} = \frac{m_{ij}}{\sum_{k=1}^m m_{ik}} = [F^T]_{ij}$$

C'est bien la formule de Markov

(Markov a donné le MLE

et $\hat{P}_T \xrightarrow{T \rightarrow \infty} P$ en vertu

des thms généraux de CV du MLE

□

)