

 ÉCOLE UNIVERSITAIRE DE PHYSIQUE ET D'INGÉNIERIE Université Clermont Auvergne	2024-2025 M2PFA UP & NP IMMAP DU Data Scientist
Statistics homework October 2024	

A propos de l'échantillonnage aléatoire

Il est parfois nécessaire de générer un échantillon de valeurs aléatoires qui suit une distribution expérimentale, qui peut posséder (ou non) une définition analytique explicite. Le but de cet exercice est d'explorer quelques-unes des techniques permettant de produire de tels échantillons.

I - Échantillonner la réciproque de la fonction cumulative.

Supposons que la fonction de répartition (ou fonction cumulative, CDF) de la distribution ciblée, \mathcal{F}_x , soit connue, ou tout du moins qu'elle puisse être dérivée numériquement des données expérimentales. Afin de générer un échantillon aléatoire de variables indépendantes X , identiquement distribuée suivant la densité de cumulative \mathcal{F}_x , on peut utiliser une transformation de la distribution continue *uniforme* dans l'intervalle $[0, 1]$, $U \sim \mathbb{I}_{[0,1]}$. Cette dernière a la propriété de coïncider avec sa propre réciproque : $\mathbf{P}(U < p) = \mathcal{F}_U(p) = p = \mathcal{F}_U^{-1}(p)$. La *méthode d'échantillonnage inverse* repose sur le théorème de la réciproque qui stipule que pour toute fonction réciproque $Q_x(p) = \mathcal{F}_x^{-1}(p)$, la variable aléatoire construite comme $X = Q_x(U)$ admet $\mathcal{F}_x(x)$ comme fonction de répartition.

Dans le cadre de cet exercice, considérons par exemple la fonction *sigmoïde* définie par:

$$\mathcal{F}_x(x; \lambda) = \frac{1}{1 + e^{-\lambda x}}.$$

avec pour paramètre $\lambda \in \mathbb{R}^+$. Cette fonction monotone croissante, $\mathcal{F}_x : x \in \mathbb{R} \rightarrow \mathcal{F}_x(x, \lambda) \in [0, 1]$ constitue une fonction de répartition valide et normalisée pour la variable aléatoire continue X . Il est facile de déterminer la fonction réciproque (ou fonction *quantile*, *percent point fonction*, PPF):

$$\mathcal{F}_x(x, \lambda) = \mathbf{P}(X < x) = p \implies x = Q_x(p, \lambda) = \mathcal{F}_x^{-1}(p, \lambda) = \frac{1}{\lambda} \ln \left(\frac{p}{1-p} \right).$$

- 1- Générez un échantillon de $N = 10^5$ probabilités aléatoires p_k , uniformément distribuées dans l'intervalle $[0, 1]$.¹ Pour chaque valeur p_k , évaluez la valeur correspondante du quantile de \mathcal{F}_x

$$x_k = \mathcal{F}_x^{-1}(p_k; \lambda = 1) = \ln(p_k / (1 - p_k)).$$

Affichez l'histogramme de densité obtenu pour l'échantillon $\{x_k\}$.

Vous avez maintenant produit un échantillon aléatoires de valeurs distribuées suivant la loi $X \sim f_x(x; \lambda = 1)$, pour la valeur de paramètre $\lambda = 1$.

- 2- Il est assez facile de dériver explicitement l'expression de la *fonction de densité de probabilité* (PDF):

$$f_x(x, \lambda) = \frac{\partial \mathcal{F}_x(x; \lambda)}{\partial x} = \frac{\lambda}{4} \operatorname{sech}^2 \left(\frac{\lambda x}{2} \right),$$

où $\operatorname{sech}(z) = 1 / \cosh(z)$ est la fonction *sécante hyperbolique*.

Superposez la courbe de densité, $f_x(x, 1)$, sur l'histogramme. Correspond-elle ?

¹ Adaptez N à vos ressources informatiques. Vous pouvez commencer avec une valeur plus faible pour développer votre code, $N=10^4$ par exemple, puis augmenter progressivement tant que le temps d'exécution reste raisonnable.

- 3- Afin de valider la distribution générée, estimons le paramètre λ de la PDF théorique en utilisant la méthode du Maximum de Vraisemblance (ML). Définissez la fonction négative du log-vraisemblance en fonction du paramètre λ , évaluée sur l'échantillon $\{x_k\}$:

$$-\ln \mathcal{L}(\lambda; \{x_k\}) = - \sum_{k=1}^N \ln [f_x(x_k; \lambda)],$$

et trouvez la valeur de l'estimateur $\hat{\lambda}$ qui minimise cette fonction.

- 4- Évaluez l'incertitude associée, $\hat{\sigma}_\lambda$, en utilisant la méthode graphique, $\Delta \ln \mathcal{L}(\lambda_\pm) = 1/2$. Est-ce que la valeur mesurée sur l'échantillon, $\hat{\lambda} \pm \hat{\sigma}_\lambda$, est compatible avec la valeur théorique attendue $\lambda = 1$?
- 5- Si oui, félicitations ! Vous savez maintenant générer un échantillon aléatoire à partir de n'importe quelle fonction de répartition.² Vous avez cependant travaillé pour rien, sinon la gloire, car cette loi de probabilité bien connue est déjà implémentée dans la librairie *scipy* sous le nom de *loi logistique*, avec toute la machinerie nécessaire à la génération aléatoire d'événements, ou d'évaluation des fonctions PDF, CDF, PPF, ou log-Likelihood ...
- Produisez un échantillon alternatif de N variables *logistiques* avec le générateur intégré (`scipy.stats.logistic.rvs(size=N)`). Comparez cet échantillon avec votre propre échantillonnage, en produisant le *graphe des résidus*, c'est-à-dire l'histogramme de la différence pour chaque *bin* des deux histogrammes, divisées par l'incertitude associée. Vous pourrez faire l'approximation, $\sigma_{n_i} \simeq \sqrt{n_i}$, pour l'incertitude sur le comptage d'événements du *bin* i (attention à la normalisation en densité). Conclure.
- 6- La loi *logistique* possède un profil symétrique avec des queues plus étalées que la loi Gaussienne-Normale (ce genre de profil est appelé *leptokurtique*). La Valeur attendue et la Variance d'une variable aléatoire *logistique* centrée, $x \sim f_x(x, \lambda)$, sont:

$$\mathbb{E}[x] = 0 \quad \text{et} \quad \mathbb{V}[x] = \frac{\pi^2}{3\lambda^2}.$$

Essayons de jouer avec cette Variance pour évaluer le nombre π .

Utilisez l'échantillon $\{x_k\}$ initialement généré avec le paramètre $\lambda = 1$, pour construire un nouvel échantillon $\{z_k\}$, où la variable aléatoire Z est définie comme la valeur centrée réduite de x multipliée par le facteur $1/\sqrt{3}$, i.e.

$$z_k = \frac{x_k - \bar{\mu}_x}{\sqrt{3\bar{\sigma}_x^2}}.$$

Les grandeurs $\bar{\mu}_x$ et $\bar{\sigma}_x^2$, sont respectivement la moyenne et la variance de l'échantillon $\{x_k\}$.³ Dans la limite des grands N , la variance de l'échantillon tend vers la variance de la distribution, $\bar{\sigma}_x^2 \simeq \mathbb{V}[x]$. La Variance de la variable aléatoire modifiée Z vaut alors $\mathbb{V}[Z] \simeq 1/3$ et sa densité de probabilité est une loi *logistique*, $f_z(z; \lambda_z)$ de paramètre $\lambda_z = \pi$.

Répétez les questions I-3 et I-4 pour évaluer l'estimateur $\hat{\lambda}_z$ qui minimise la fonction négative du log-vraisemblance, $-\ln \mathcal{L}(\lambda; \{z_k\})$, évaluée sur l'échantillon modifié $\{z_k\}$, ainsi que l'incertitude associée.⁴ Est-ce que la mesure est compatible avec la valeur connue de π ? (et vous pouvez ainsi conclure que les méthodes d'échantillonnage sont généralement très peu efficaces pour déterminer un nombre significatif de décimales de π)

²Dans les cas où l'expression analytique de la fonction quantile ne peut pas être dérivée facilement, il est toujours possible d'interpoler la distribution expérimentale pour déterminer les valeurs de x correspondant à n'importe quelle valeur $p = \mathcal{F}_x(x)$ de la fonction.

³Vous pouvez rapidement obtenir ces valeurs en important le module python *statistics* qui fournit: $\bar{\mu}_x = \text{statistics.mean}(X)$ et $\bar{\sigma}_x^2 = \text{statistics.variance}(X)$.

⁴L'incertitude sur la détermination de $\bar{\sigma}_x^2$ est négligée ici.

II - La méthode d'échantillonnage accept-reject.

La méthode accept-reject est une technique *Monte-Carlo* qui permet de “sculpter” un échantillon aléatoire d'une densité cible, $g_x(x)$, à partir d'une densité de référence, $f_x(x)$, qui vérifie le critère

$$\alpha f_x(x) \geq g_x(x) \quad \forall x \in \text{supp}(g_x),$$

où α est un facteur d'échelle constant qui assure l'inégalité en tout x .

La méthode consiste à

- i) générer une probabilité aléatoire p uniformément distribuée dans l'intervalle $[0, 1]$.
- ii) générer une valeur aléatoire x qui suit la densité de référence $f_x(x)$.
- iii) si le rapport des densités $g_x(x)/\alpha f_x(x) \geq p$ la valeur x est acceptée (rejetée sinon).
- iv) répéter cette procédure N fois.

La méthode accept-reject nécessite généralement d'importants échantillons de référence, en particulier pour peupler les régions de faible densité de la distribution cible. Plus les profils entre distribution de référence et cible sont proches, plus le taux de rejection est faible.

Essayons de générer un échantillon aléatoire distribué suivant la loi *Normale*, en utilisant la loi *logistique* plus étalée comme référence.

- 1- Le maximum de densité vaut $f_x(0; 1) = 1/4$ pour la loi *logistique*, et $\mathcal{N}_x(0; 0, 1) = 1/\sqrt{2\pi}$ pour la loi *Normale*. Nous pouvons choisir la valeur minimale du facteur d'échelle $\alpha_0 = 4/\sqrt{2\pi}$ telle que la distribution *logistique* “englobe” la loi de densité *Normale* pour tout x .
Tracer les deux fonctions, $\alpha_0 f_x(x, 1)$ et $\mathcal{N}_x(x; 0, 1)$, sur un même graphe pour vérifier que l'on a bien $\alpha_0 f_x(x) \geq \mathcal{N}_x(x) \quad \forall x$.⁵
- 2- Pour chacune des valeurs de la distribution *logistique* de l'échantillon $\{x_k\}$, ($k = 1, N$)
 - (i) évaluez la rapport des densité $r(x_k) = \mathcal{N}_x(x_k)/\alpha_0 f_x(x_k, 1)$.
 - (ii) générez une probabilité aléatoire p uniformément distribuée dans l'intervalle $[0, 1]$
 - (iii) si $r(x_k) \geq p$ gardez x_k dans le sous-échantillon des valeurs acceptées $\{x_m\}_a$, $m = 1, M \leq N$.
Quel est le taux d'acceptance M/N ?
- 3- Construisez l'histogramme de densité pour le sous-échantillon $\{x_m\}_a$. Vérifiez le profil *Normal* en réalisant l'ajustement d'une loi gaussienne sur le sous-échantillon (avec `scipy.stats.norm.fit`, par exemple). Donner les valeurs obtenues des paramètres ajustés $\hat{\mu}$ et $\hat{\sigma}$ et tracez la courbe de densité, $\mathcal{N}(x; \hat{\mu}, \hat{\sigma})$, sur l'histogramme.
- 4- Bien que généralement un choix sous-optimal, la loi *Uniforme* dans un intervalle fini $[a, b]$, englobant tout autre distribution avec le bon facteur d'échelle, peut être utilisée comme distribution de référence pour produire n'importe quelle distribution aléatoire sur un intervalle fini. Déterminez le facteur d'échelle minimal α_0 et produisez un échantillon de valeurs aléatoires sur l'intervalle $[-10, 10]$ pour la loi *logistique* $x \sim f_x(x; 1)$, à partir d'un échantillon de N valeurs aléatoires uniformément distribuées. Quel est le taux d'acceptance M/N ? Essayez de trouver un meilleur choix pour la distribution de référence que la loi uniforme, pour produire un échantillon *logistique*.

⁵La méthode fonctionne identiquement pour tout facteur d'échelle $\alpha \geq \alpha_0$, mais plus la valeur est grande plus le taux de rejection augmente.