

Echantillonnage Décisionnel - README

abelmasson

February 2025

1 Introduction

L'échantillonnage est une méthode consistant à sélectionner un sous-ensemble d'observations (= un échantillon) représentatif d'une population plus large pour répondre à une question. Selon cette question, différents critères peuvent guider la construction de l'échantillon. Dans la plupart des cas, un échantillon de grande taille est préférable, parce qu'il améliore la représentativité de la population et la précision des réponses. Toutefois, dans certaines situations, notamment lorsque l'acquisition de chaque observation est coûteuse, on cherchera plutôt à limiter la taille de l'échantillon. L'échantillonnage décisionnel (ou séquentiel) est particulièrement adapté à ces situations en ce qu'il repose sur une collecte progressive des observations, une par une, afin d'obtenir un échantillon de taille strictement suffisante pour fournir une réponse fiable.

Exemple d'application à l'infestation d'une serre par un insecte

On cherche à surveiller le risque d'infestation d'une serre de fraises en constituant un échantillon de plants de fraisier. Chaque plant est consciencieusement examiné par un technicien pour déterminer la présence ou l'absence de l'insecte. Dans ce contexte, un faible nombre de détection peut être le signal d'une infestation, alors que l'insecte recherché est généralement rare et souvent absent. Ainsi, le nombre de plants à examiner avant de pouvoir statuer sur l'infestation de la serre peut rapidement devenir ingérable. Ce script propose une méthode d'échantillonnage décisionnel visant à optimiser, c'est-à-dire à réduire au maximum le nombre de plants à examiner.

1.1 Formulation

1.1.1 Formulation du modèle

Pour chaque lot de plants (i.e. chaque serre), on cherche à déterminer si la fréquence de l'insecte dépasse un certain seuil au delà duquel on peut considérer que la serre est infestée. On note p la probabilité de retrouver un insecte sur un plant, et p^* le seuil critique à ne pas dépasser (par exemple 1 % des plants soit $p^* = 0.01$). On suppose que p suit une loi *Beta* et que X le nombre de

plants infestés suit une loi Binomiale de paramètres N le nombre total de plants examinés et p . Ainsi :

$$X \sim \text{Binomial}(N, p)$$

avec

$$p \sim \text{Beta}(\alpha, \beta)$$

1.1.2 Formulation du test

On pose **(H)** l'hypothèse : "la probabilité p de trouver un insecte sur un plant dépasse le seuil critique p^* "

On cherche à tester **(H)** en fonction des résultats des examens de plants, c'est-à-dire en fonction des observations de X . Pour ce faire on se munit d'un seuil de confiance e permettant d'accepter l'hypothèse **(H)** (lorsque $P(Hvraie) > 1 - e$), de la rejeter (lorsque $P(Hvraie) < e$), ou de rester indécis si les résultats des examens ne permettent pas de prendre une décision fiable (lorsque $1 - e > P(Hvraie) > e$).

Rappel sur la loi Beta

La loi *Beta* est une distribution continue définie par deux paramètres α et β . Elle est classiquement utilisée pour représenter des proportions ou des probabilités. En particulier lorsque l'on s'intéresse à des données distribuées selon une loi binomiale ($X \sim \text{Binomial}(N, p)$), la loi *Beta* est adaptée à la représentation de la probabilité p pour plusieurs raisons :

- La loi *Beta* est définie sur l'intervalle $[0, 1]$.
- La loi *Beta*(1, 1) est identique à la loi uniforme sur $[0, 1]$ et constitue une *prior* adaptée lorsqu'aucune autre information sur p n'est disponible.
- Lorsque des observations de X sont disponibles, il existe une expression analytique pour en informer la distribution de p de façon simple et intuitive : La *posterior* de p sera la loi *Beta*($1 + \alpha$, $1 + \beta$) avec α et β les nombres de succès et d'échecs observés. On dit que la loi *Beta* est la conjuguée de la loi Binomiale.

1.2 Protocole d'échantillonnage décisionnel

Comme expliqué en introduction, le protocole d'échantillonnage décisionnel repose sur l'intégration des observations une par une. Plus précisément il reestime la distribution de probabilité de p et reitere le test sur **(H)** à chaque nouvel examen de plant. Le détail du protocole est le suivant :

- On part sans information a priori sur la distribution de p , soit $p \sim \text{Beta}(1, 1)$

- A chaque nouvel examen de plant :
 - On actualise la distribution de p , en ajoutant 1 à α si le test est positif, 1 à β s'il est négatif (voir rappel sur la loi Beta).
 - On calcule la probabilité $P(H \text{ vraie})$ que p dépasse le seuil critique p^* .
 - * Si $P(H \text{ vraie})$ est supérieure à $1 - e$ alors on conclut que **le lot est infesté** et on s'arrête.
 - * Si $P(H \text{ vraie})$ est inférieure à e on conclut que **le lot n'est pas infesté** et on s'arrête.
 - * Si $P(H \text{ vraie})$ est comprise entre e et $1 - e$ enfin, un **examen supplémentaire est nécessaire**, on recommence la boucle avec un nouvel examen.
- Si l'on atteint le nombre de plants dans le lot avant d'avoir pu conclure que le lot est ou n'est pas infesté, on conclut d'une **indécision**.

2 Expérience Interactive

Le code .Rmd présenté dans ce *reposiroty* contient un RShiny pour implémenter et visualiser le protocole d'échantillonnage décisionnel. Ci-après une capture d'écran de l'interface d'expérience interactive (Figure 1). Elle en donne une représentation graphique comprenant :

- Un graphique avec en ordonnées le nombre de plant examinés infestés (N_+) et en abscisse le nombre de plants examinés sains (N_-). Sur ce graphique on représente en rouge la zone où $P(H \text{ vraie} | X = \frac{N_+}{N_+ + N_-}) > 1 - e$ et en bleu $P(H \text{ vraie} | X = \frac{N_+}{N_+ + N_-}) < e$.
- Un graphique donnant la distribution de probabilité p qui correspond à une distribution $Beta(1 + N_+, 1 + N_-)$.

L'expérience interactive prend automatiquement fin dès qu'une décision a été atteinte.

