

# Modèles discrets de distributions fréquents

Antonio Falcó

- 1 Variables aléatoires
- 2 Distribution Bernoulli
- 3 Espérance mathématique
- 4 Échantillonnage aléatoire simple
- 5 Distribution Binomiale
  - Application pratique
  - Estimateurs
- 6 Distribution de Poisson
  - Application pratique

## Hypothèse

- 1 Soit  $\Omega$  l'ensemble des individus associé à une expérience scientifique et
- 2 Pr la loi de probabilité associé à cette expérience:

$$\text{Pr} : \text{Événements dans } \Omega : \longrightarrow [0, 1].$$

## La mesure quantitative

Soit  $X : \Omega \longrightarrow \mathbb{R}$  une grandeur qu'on utilise pour étudier les individus de  $\Omega$  et qu'on appelle **variable aléatoire** si pour chaque pair des mesures  $x, x' \in X(\Omega)$  avec  $x \leq x'$ , qu'on peut observer dans quelque expérience réalisé à  $\Omega$ , l'ensemble

$$\{\omega \in \Omega : x < X(\omega) \leq x'\}$$

est un événement à  $\Omega$ .

## Exemple: Taille

Soit  $\Omega$  une population d'individus, et

$$\text{Taille} : \Omega \longrightarrow \mathbb{R}$$

la mesure de la taille en cm sur chaque individu dans la population.

- Soit  $x' = 0$  et  $x = -120$ , alors

$$\{\omega \in \Omega : -120 < \text{Taille}(\omega) \leq 0\} = \emptyset,$$

- Soit  $x = 0$  et  $x' = 160.5$ , alors

$$\{\omega \in \Omega : 0 < \text{Taille}(\omega) \leq 160.5\}.$$

Est-ce qu'on peut calculer

$$\Pr(\{\omega \in \Omega : 0 < \text{Taille}(\omega) \leq 160.5\})?$$

## Exemple; L'âge

Soit  $\Omega$  une population d'individus, et

$$\text{Age} : \Omega \longrightarrow \mathbb{R}$$

la mesure de l'âge en années (nombres entiers non négatives) sur chaque individu dans la population.

- Soit  $x' = 0$  et  $x = -12$ , alors

$$\{\omega \in \Omega : -12 < \text{Age}(\omega) \leq 0\} = \emptyset,$$

- Soit  $x = 0$  et  $x' = 16$ , alors

$$\{\omega \in \Omega : 0 < \text{Age}(\omega) \leq 16\}.$$

Est-ce qu'on peut calculer

$$\Pr(\{\omega \in \Omega : 0 < \text{Age}(\omega) \leq 160\})?$$

## Caractéristiques

- 1 On travaille avec des variables quantitatives,
- 2 Taille( $\Omega$ )  $\subset [0, 1000]$  cm, alors il est une variable continue.
- 3 Age( $\Omega$ )  $\subset \{1, 2, \dots, 1000\}$  années, alors il est une variable discrète.

## Conséquence

- 1 Si  $X$  est une variable aléatoire continue on peut calculer

$$\Pr(\{\omega \in \Omega : x < X(\omega) \leq x'\}) \equiv \Pr(x < X \leq x')$$

pour tout  $x, x' \in \mathbb{R}$ .

- 2 Si  $X$  est une variable aléatoire discrete on peut calculer

$$\Pr(\{\omega \in \Omega : X(\omega) = k\}) \equiv \Pr(X = k)$$

pour tout  $k, \in \mathbb{Z} = \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$ .

# Distribution Bernoulli

Considérons une expérience qui n'a que deux issues possibles (ex: jet d'une pièce, on étudie la prévalence d'une maladie dans une population), et convenons d'appeler  $S$  la première issue ("succès") et  $\bar{S}$  la seconde ("échec"). Soit

$$X(\omega) := \begin{cases} 1 & \text{si } \omega \in S \text{ "succès"} \\ 0 & \text{si } \omega \in \bar{S} \text{ "échec"} \end{cases}$$

Les issues possibles dans la population sont  $\{S\}$  et  $\{\bar{S}\}$  :  $S \cap \bar{S} = \emptyset$  sont incompatibles et  $S \cup \bar{S} = \Omega \equiv$  population . La distribution des probabilités est:

$$\begin{aligned} \Pr(X = 1) &= \Pr(S) = p \\ \Pr(X = 0) &= \Pr(\bar{S}) = 1 - p \end{aligned}$$

L'ensemble des modalités est  $X(\Omega) = \{0, 1\}$ .

## Variable aléatoire discrete

Soit  $X : \Omega \longrightarrow \mathbb{R}$  une v.a. discrete i.e.  $X(\Omega)$  est un ensemble dénombrable.

- ❶ On appelle  $f$  la fonction de  $f : X(\Omega) \longrightarrow [0, 1]$  donné par

$$f(k) := \Pr(X = k)$$

fonction de densité de probabilité.

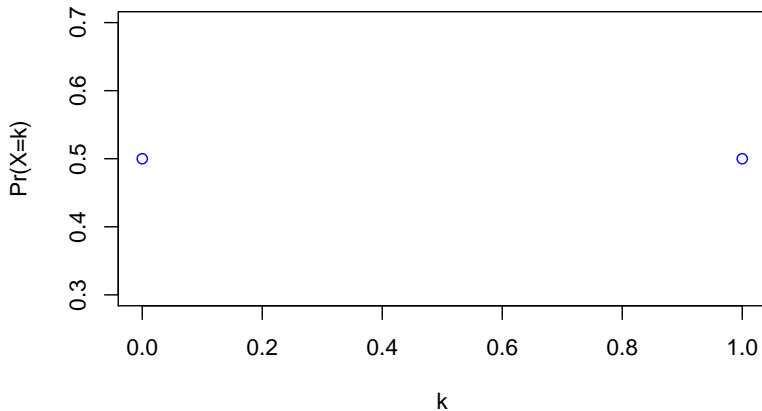
- ❷ On a la propriété

$$\sum_{k \in X(\Omega)} f(k) = \sum_{k \in X(\Omega)} \Pr(X = k) = 1.$$

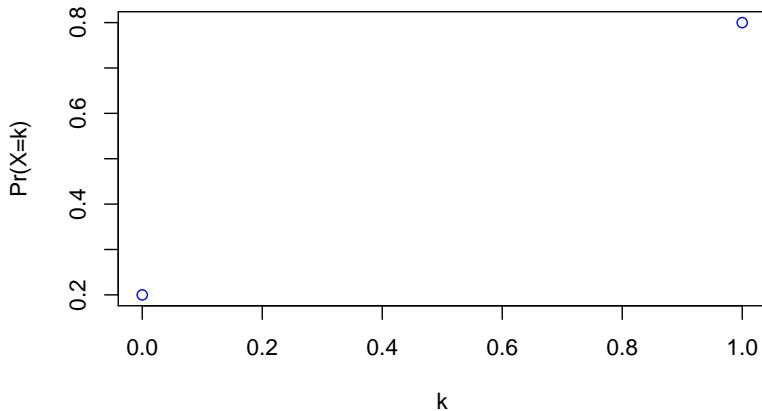
- ❸  $\Pr(X = k)$  joue le rôle de la fréquence relative pour la modalité  $k$ .



## Distribution de Bernoulli $p=0.5$



## Distribution de Bernoulli $p=0.8$



## Fréquence cumulée

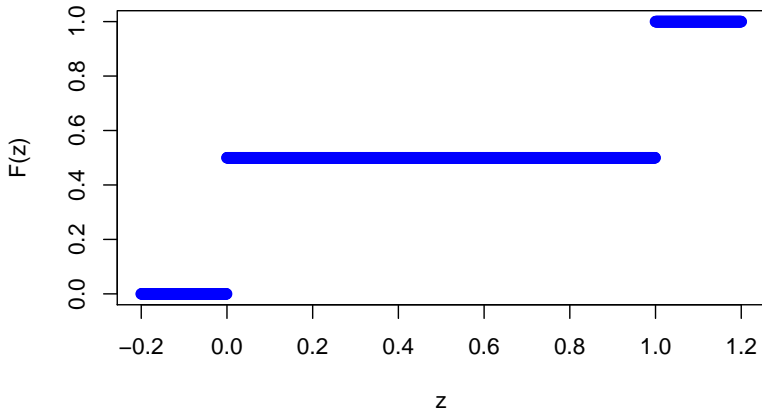
On peut construire la fonction de distribution de  $X$  en utilisant:

$$F : \mathbb{R} \longrightarrow [0, 1]$$

définie par

$$F(z) = \sum_{\substack{k \leq z \\ k \in X(\Omega)}} \Pr(X = k).$$

## Fonction de distribution $p=0.5$



## L'espérance mathématique d'une v.a. discrète

Soit  $X$  un v.a. discrète et  $g$  une fonction n'importe lequel. Alors on va définir l'espérance mathématique de  $g(X)$  comme:

$$E[g(X)] = \sum_{k \in X(\Omega)} g(k) \cdot \Pr(X = k)$$

Moyenne:  $g(X) = X$

On appelle moyenne de  $X$  à l'espérance mathématique de  $g(X) = X$  :

$$\mu := E[X] = \sum_{k \in X(\Omega)} k \cdot \Pr(X = k)$$

Variance:  $g(X) = (X - E[X])^2$

On appelle variance de  $X$  à l'espérance mathématique de  $g(X) = (X - E[X])^2$  :

$$\sigma^2 := E[(X - E[X])^2] = \sum_{k \in X(\Omega)} (k - E[X])^2 \cdot \Pr(X = k)$$

## Exemple: Distribution de Bernoulli

- Moyenne:  $E[X] = p$

$$E(X) = 1 \cdot \Pr(X = 1) + 0 \cdot \Pr(X = 0) = 1 \cdot p + 0 \cdot (1 - p) = p$$

- Variance:  $\text{Var}(X) = E[(X - E[X])^2] = E[(X - p)^2] = p(1 - p)$

$$\begin{aligned} E[(X - p)^2] &= (1 - p)^2 \cdot \Pr(X = 1) + (0 - p)^2 \cdot \Pr(X = 0) \\ &= (1 - p)^2 \cdot p + p^2 \cdot (1 - p) \\ &= p(1 - p)[(1 - p) + p] \\ &= p(1 - p) \end{aligned}$$

## L'espérance mathématique est un opérateur linéaire

Soit  $X_1 = g_1(X)$ ,  $X_2 = g_2(X)$ . Alors

$$\begin{aligned} E[X_1 + X_2] &= \sum_{k \in X(\Omega)} (g_1(k) + g_2(k)) \cdot \Pr(X = k) \\ &= \sum_{k \in X(\Omega)} g_1(k) \cdot \Pr(X = k) + \sum_{k \in X(\Omega)} g_2(k) \cdot \Pr(X = k) \\ &= E[X_1] + E[X_2]. \end{aligned}$$

Soit  $\lambda \in \mathbb{R}$ , alors

$$\begin{aligned} E[\lambda g_1(X)] &= \sum_{k \in X(\Omega)} \lambda \cdot g_1(k) \cdot \Pr(X = k) = \lambda \cdot \sum_{k \in X(\Omega)} g_1(k) \cdot \Pr(X = k) \\ &= \lambda \cdot E[g_1(X)]. \end{aligned}$$



## Propriétés

- ❶ Soit  $X_1, \dots, X_n$  des v.a. avec  $X_1 = g_1(X), \dots, X_n = g_n(X)$  et  $\lambda_1, \dots, \lambda_n$  des nombres réels. Alors

$$E \left[ \sum_{i=1}^n \lambda_i \cdot X_i \right] = \sum_{i=1}^n \lambda_i \cdot E[X_i].$$

- ❷ Soit  $g$  la fonction constant égal à  $c$  i.e  $Y = g(X) = c$  alors

$$\begin{aligned} E[Y] &= E[c] = \sum_{k \in X(\Omega)} g(k) \cdot \Pr(X = k) = \sum_{k \in X(\Omega)} c \cdot \Pr(X = k) \\ &= c \cdot \overbrace{\sum_{k \in X(\Omega)} \Pr(X = k)}^{=1} = c \end{aligned}$$

- La caractéristique de la population  $X$  dans le modèle Bernoulli dépend de la distribution de population, noté par  $\mathcal{B}(p)$ , et est déterminée par la valeur de le paramètre  $p$ .
- Pour noté que la variable aléatoire  $X$  suivi une distribution Bernoulli avec paramètre  $p$  on utilise la notation:

$$X \sim \mathcal{B}(p),$$

- **L'inférence statistique** consiste en la détermination de la distribution de population  $\mathcal{F}$  ( $\mathcal{F} = \mathcal{B}(p)$  dans le cas Bernoulli) et de ses caractéristiques (moyenne, variance, quantiles, ...) à partir des observations  $x_1, \dots, x_n$  sur l'échantillon, ainsi qu'en l'étude de la précision avec laquelle ces caractéristiques sont déterminées.

# Échantillonnage aléatoire simple

- Afin d'obtenir un échantillon représentatif d'une population, il est nécessaire de le tirer de façon aléatoire.
- L'exemple classique d'échantillonnage aléatoire simple consiste à placer des billets contenant les noms de tous les individus de la population dans une urne et de tirer des billets au hasard sans remise.
- Dans la pratique, ce principe est mis en oeuvre à l'aide de logiciels permettant de générer des nombres aléatoires, sur la base desquels on sélectionne des individus à partir d'une liste (ex.: annuaire téléphonique).

# Échantillonnage aléatoire simple

- Il existe des procédés d'échantillonnage plus sophistiqués, comme par exemple l'échantillonnage aléatoire stratifié, où l'on échantillonne séparément dans des sous-populations (appelées strates), par exemple pour garantir d'avoir des proportions d'individus de chaque strate qui soient conformes aux proportions de la population.
- Dans ce cours, nous nous concentrerons sur l'échantillonnage aléatoire simple.

- Considérons un ensemble de  $n$  individus tirés d'une population à l'aide d'un échantillonnage aléatoire simple et intéressons-nous à une caractéristique  $X$  de ces individus.
- On considère les mesures de  $X$  que nous allons faire sur chaque individu comme des variables aléatoires  $X_1, \dots, X_n$  et on fait les hypothèses suivantes:
  - ① Les variables  $X_1, \dots, X_n$  sont indépendantes,
  - ② Les variables  $X_1, \dots, X_n$  ont toutes la même distribution  $\mathcal{F}$ , où  $\mathcal{F}$  est la distribution (inconnue) de la caractéristique d'intérêt dans la population.
- On résume ces deux hypothèses en disant que  $X_1, \dots, X_n$  sont indépendantes et identiquement distribuées selon  $\mathcal{F}$ , ce qu'on note

$$X_1, \dots, X_n \text{ i.i.d. } \sim \mathcal{F}$$

# Distribution Binomiale

Soit

$$X_1, \dots, X_n \text{ i.i.d. } \sim \mathcal{F} = \mathcal{B}(p)$$

et on considère:

$$Z := X_1 + \dots + X_n$$

où

$Z$  = nombre de succès parmi les  $n$  répétitions de  $X$

L'ensemble des modalités pour  $Z$  est

$$Z(\Omega) = \{0, 1, 2, \dots, n\}.$$

On dit que  $Z$  a (ou suit) une **distribution binomiale de paramètres  $n$  et  $p$**  :

$$Z \sim \mathcal{B}(n, p).$$

Pour  $n = 1$  on a la distribution  $\mathcal{B}(p) = \mathcal{B}(p, 1)$ .

# Distribution des probabilités

- La distribution de  $Z \sim \mathcal{B}(n, p)$  est donnée par

$$\Pr(Z = k | n, p) = \binom{n}{k} p^k (1 - p)^{n-k} \quad k = 0, 1, \dots, n$$

où le coefficient binomial est défini comme

$$\binom{n}{k} = \frac{n!}{(n-k)!k!}$$

et  $n!$  ( $n$  factoriel) est défini comme

$$n! = n(n-1) \cdots 2 \cdot 1.$$

Par convention  $0! = 1$ .

Avec le logiciel R si  $Z \sim \mathcal{B}(n, p)$  on peut calculer la probabilité avec la fonction de densité de probabilité: `dbinom()` :

$$\Pr(Z = k | n, p) = \text{dbinom}(k, n, p)$$

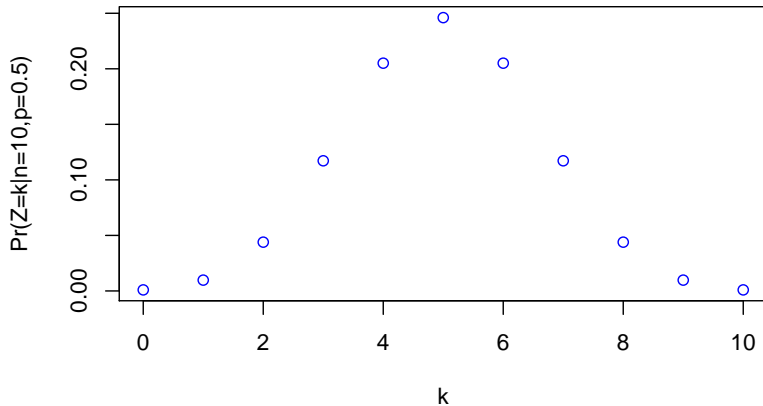
### Exemple

Soit  $n = 10$  et  $p = 0.05$ , alors

$$\Pr(Z = 1 | n = 10, p = 0.05) = \text{dbinom}(1, 10, 0.05) = 0.0104751$$



## Fonction de densité de probabilité $B(n,p)$



Avec le logiciel R si  $Z \sim \mathcal{B}(n, p)$  on peut calculer la probabilité cumulée avec la fonction de distribution `pbinom()` :

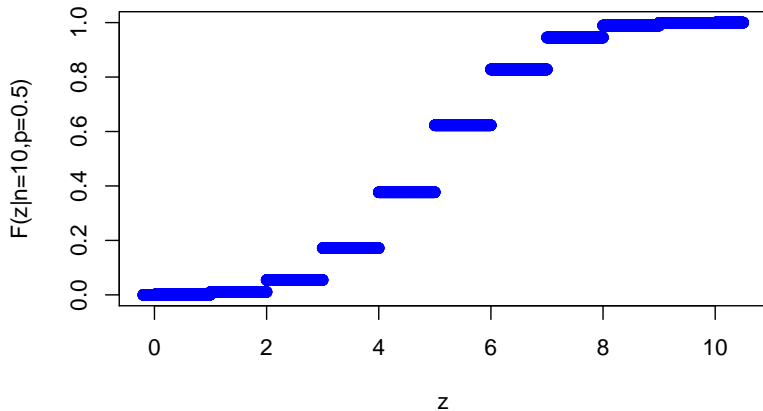
$$F(z|n, p) = \Pr(Z \leq z|n, p) = \text{pbinom}(z, n, p)$$

### Exemple

Soit  $n = 10$  et  $p = 0.05$ , alors

$$\begin{aligned} F(3.5|n = 10, p = 0.05) &= \Pr(Z \leq 3.5|n = 10, p = 0.05) \\ &= \text{pbinom}(3.5, 10, 0.05) = 0.171875 \end{aligned}$$

## Fonction de distribution



## Application pratique

- On sait que 10% des personnes montrant un certain symptôme sont porteuses d'une maladie donnée. Dans la ville considérée cette maladie peut être qualifiée d'épidémie.
- Le diagnostic précis dépend d'un test sanguin qui est malheureusement assez coûteux.
- En conséquence, les hématologues attendent d'avoir reçu la visite de  $n$  porteurs du symptôme avant de réaliser ce test avec un mélange du sang de ces  $n$  personnes.
- Si aucune d'entre elles n'est malade, le test est négatif. Dans le cas contraire, le test est positif et le médecin doit effectuer des tests individuels sur le sang de chacun des patients pour détecter qui est malade.

# Application pratique

## Données

- Population: L'ensemble des personnes montrant un certain symptôme.
- $S$  = sont porteuses d'une maladie donnée.
- $X = 1$  si la personne est porteuse et  $X = 0$  au contraire.
- $\Pr(X = 1) = \Pr(S) = p = 0.1$ .
- $Z = X_1 + \dots + X_n$  nombre de porteuses de la maladie dans un échantillon de  $n$ -personnes avec le symptôme.
- $Z \sim \mathcal{B}(n, p = 0.1)$ .

## Application pratique

### Question

1. Quelle est la probabilité que le test collectif soit négatif lorsque  $n = 2$ ,  $n = 4$ ,  $n = 6$  et  $n = 10$ ?

Si

$$Z = 0$$

alors le test est négatif. La probabilité que le test collectif soit négatif:

$$\Pr(Z = 0 | n = 2, p = 0.01) = \text{dbinom}(0, 2, 0.1) = 0.81$$

$$\Pr(Z = 0 | n = 4, p = 0.01) = \text{dbinom}(0, 4, 0.1) = 0.6561$$

$$\Pr(Z = 0 | n = 6, p = 0.01) = \text{dbinom}(0, 6, 0.1) = 0.531441$$

$$\Pr(Z = 0 | n = 10, p = 0.01) = \text{dbinom}(0, 10, 0.1) = 0.4304672$$

## Application pratique

### Question

2. Soit  $Y$  le nombre de tests à effectuer en général pour un groupe de  $n$  patients. Donner la loi de  $Y$  pour  $n = 2$ ,  $n = 4$ ,  $n = 6$  et  $n = 10$ .

Soit

$$Y = \begin{cases} 1 & \text{si } Z_n = 0 \\ n + 1 & \text{si } Z_n \neq 0. \end{cases}$$

Alors, l'ensemble des modalités de  $Y$  est  $\{1, n + 1\}$  et la distribution des probabilités:

$$\Pr(Y = 1) = \Pr(Z = 0|n, p) = p,$$

et

$$\Pr(Y = n + 1) = \Pr(Z > 0|n, p) = 1 - \Pr(Z = 0|n, p) = 1 - p.$$

La variable  $Y \sim \mathcal{B}(p = \Pr(Z = 0|n, p))$ .

## Moyenne et variance $Z \sim \mathcal{B}(n, p)$

- Moyenne:  $E[Z] = np$

$$\begin{aligned} E[Z] &= E[X_1 + \cdots + X_n] = E[X_1] + \cdots + E[X_n] \\ &= \overbrace{p + \cdots + p}^{n\text{-fois}} = np. \end{aligned}$$

- Variance:  $\text{Var}(Z) = E[(Z - E[Z])^2] = np(1 - p)$ .

$$\begin{aligned} \text{Var}(Z) &= \text{Var}(X_1 + \cdots + X_n) = \text{Var}(X_1) + \cdots + \text{Var}(X_n) \\ &= \overbrace{p(1 - p) + \cdots + p(1 - p)}^{n\text{-fois}} = np(1 - p). \end{aligned}$$



# L'estimateur de la moyenne

Soit  $X \sim \mathcal{F}$ , telle que  $\mu = E[X]$  et  $\sigma^2 = \text{Var}(X)$ , et soit

$$X_1, \dots, X_n \text{ i.i.d. } \sim \mathcal{F}.$$

Un **estimateur** pour la moyenne  $\mu = E[X]$  est la **moyenne empirique**:

$$\bar{X} = \frac{X_1 + \dots + X_n}{n}.$$

## Definition

On appelle **estimateur** a toute fonction d'un échantillonnage aléatoire simple:

$$\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$$

## Propriétés de la moyenne empirique

①  $E[\bar{X}] = \mu.$

$$E[\bar{X}] = E\left[\frac{X_1 + \cdots + X_n}{n}\right] = \frac{1}{n}E[X_1 + \cdots + X_n] = \frac{1}{n}n\mu = \mu.$$

②  $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}.$

$$\begin{aligned}\text{Var}(\bar{X}) &= \text{Var}\left(\frac{X_1 + \cdots + X_n}{n}\right) = \frac{1}{n^2}\text{Var}(X_1 + \cdots + X_n) \\ &= \frac{1}{n^2}n\sigma^2 = \frac{\sigma^2}{n}.\end{aligned}$$

Si  $n \rightarrow \infty$  alors  $\text{Var}(\bar{X}) \rightarrow 0$ , et une conséquence  $E[\bar{X}] \rightarrow \mu$ .

Une variable  $X$  suit une distribution de Poisson de paramètre  $\lambda > 0$  ce qu'on note  $X \sim \mathcal{P}(\lambda)$  si

$$\Pr(X = k|\lambda) = e^{-\lambda} \frac{\lambda^k}{k!} \quad k = 0, 1, 2, \dots$$

Les modalités d'une variable Poisson sont donc tous les entiers positifs plus 0. Avec le logiciel R on utilise

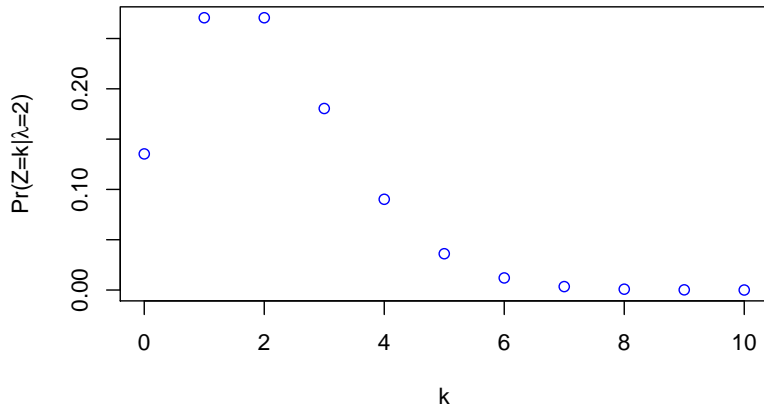
$$\Pr(X = k|\lambda) = \text{dpois}(k, \lambda),$$

et

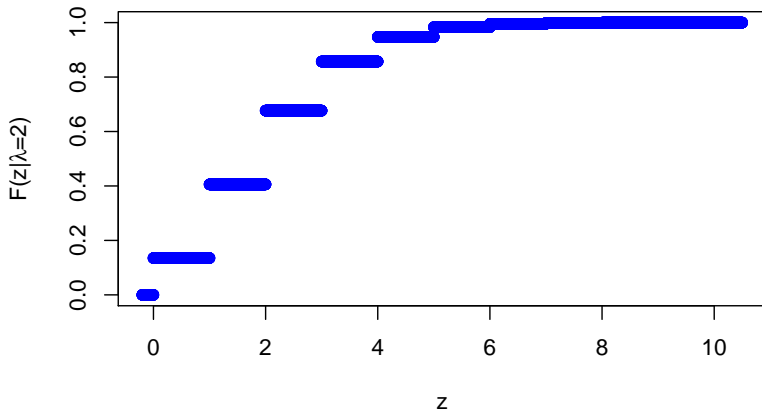
$$F(z|\lambda) = \Pr(X \leq z|\lambda) = \text{ppois}(z, \lambda).$$

$$X \sim \mathcal{P}(\lambda = 2) :$$

## Fonction de densité de probabilité



## Fonction de distribution



## Propriété

Soit  $X \sim \mathcal{P}(\lambda)$ . Alors

$$\Pr(X = 0) = e^{-\lambda} \frac{\lambda^0}{0!} = e^{-\lambda}.$$

En conséquence,

$$\lambda = -\ln(\Pr(X = 0|\lambda)) = \ln\left(\frac{1}{\Pr(X = 0|\lambda)}\right)$$

- Dans la pratique la distribution de Poisson est souvent utilisée pour modéliser des données de comptage, par exemple le nombre de nouveaux cas de cancer dans une certaine région pendant une certaine période de temps (en épidémiologie on appelle ce nombre *l'incidence*).

- La loi de Poisson est souvent utilisée pour approximer certaines lois discrètes.
- On l'appelle aussi loi des événements rares.
- En effet, si  $X$  est le nombre de fois où apparaît un événement de probabilité très petite ( $p$ ), alors la loi de  $X$  peut être approximée par une loi de poisson.
- Si  $n$  est grand et  $p$  petite alors

$$\mathcal{B}(n, p) \sim \mathcal{P}(\lambda = np).$$

C'est-à-dire, soit  $X \sim \mathcal{B}(n, p)$  et  $\hat{X} \sim \mathcal{P}(\lambda = np)$  alors

$$\Pr(X = k | n, p) \approx \Pr(\hat{X} = k | \lambda = np).$$



### Exemple

Soit  $n = 100$  et  $p = 0.06$ , en conséquence  $\lambda = n \cdot p = 6$  :

$$\Pr(X = 3 | n = 100, p = 0.06) = \text{dbinom}(3, 100, 0.06) = 0.0864103$$

et

$$\Pr(\hat{X} = 3 | \lambda = 6) = \text{dpois}(3, 6) = 0.0892351$$

### Remarque

Les conditions usuelles sous lesquelles on considère que la qualité de l'approximation est raisonnable sont les suivantes:  $n > 30$ , et  $np > 5$ .

# Applications de la loi de Poisson

- Le nombre de coquilles par page ou groupe de pages d'un livre,
- le nombre d'individus dépassant l'âge de 100 dans une communauté humaine,
- le nombre de faux numéros téléphoniques composés en un jour,
- le nombre de paquets de biscuits pour chien vendus dans un magasin donné en l'espace d'un jour,
- le nombre de clients pénétrant dans un bureau de poste donné en l'espace d'un jour,
- le nombre de décès attribués à la fièvre typhoïde sur une longue période, par exemple un an.

## Moyenne et variance $X \sim \mathcal{P}(\lambda)$

- Moyenne:  $E[X] = \lambda$ .

$$\begin{aligned} E[X] &= \sum_{k=0}^{\infty} k \Pr(X = k) = \sum_{k=0}^{\infty} k \left( e^{-\lambda} \frac{\lambda^k}{k!} \right) = e^{-\lambda} \sum_{k=0}^{\infty} k \frac{\lambda^k}{k!} \\ &= e^{-\lambda} \sum_{k=1}^{\infty} k \frac{\lambda^k}{k!} = e^{-\lambda} \sum_{k=1}^{\infty} \frac{\lambda^k}{(k-1)!} = e^{-\lambda} \lambda \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} \\ &= e^{-\lambda} \lambda \sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{-\lambda} \lambda e^{\lambda} = \lambda. \end{aligned}$$

- Variance:  $\text{Var}(X) = \lambda$ .

# Application pratique

## Interprétation

Si  $\lambda = E[X]$  et le nombre d'événements espérés par unité du temps  $t$  alors  $\frac{\lambda}{t}$  est le nombre d'événements espérés dans le période du temps  $t$ .

Supposons que le nombre de décès dus à la fièvre typhoïde sur une période d'un an correspond à une distribution de Poisson avec le paramètre  $\lambda = 4.5$ . Quelle est la distribution de probabilité du nombre de décès sur une période de 6 mois? Une période de 3 mois?

## Application pratique

Soit

$X$  = nombre de décès dus à la fièvre typhoïde sur une période de 6 mois

Si le paramètre pour un an est

$$\lambda = 4.5 \times 1 \text{ an}$$

alors le paramètre pour 6 mois est

$$\lambda = 4.5 \times \frac{1}{2} \text{ an} = 2.25 \times 6 \text{ mois} .$$

Alors  $X \sim \mathcal{P}(\lambda = 2.25)$ .

# Application pratique

## Distribution de Poisson

