

# RFIDEC — cours 2 : Échantillons, estimations ponctuelles

Christophe Gonzales

LIP6 – Université Paris 6, France

## Plan du cours n°2

- 1 Lois des grands nombres
- 2 Théorème central-limite
- 3 Estimation ponctuelle à partir d'échantillons
- 4 Biais dans les estimations

RFIDEC — cours 2 : Échantillons, estimations ponctuelles

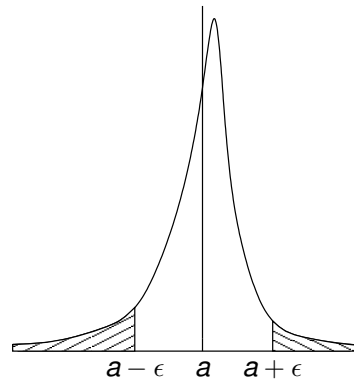
2/37

## Convergence en probabilité

### Définition

- $(X_n)_{n \in \mathbb{N}}$  : suite de variables
- $a$  : constante
- $(X_n)$  **converge en probabilité** vers  $a$  si, pour tout  $\epsilon > 0$  la probabilité que l'écart absolu entre  $X_n$  et  $a$  dépasse  $\epsilon$  tend vers 0 quand  $n \rightarrow \infty$  :

$$\lim_{n \rightarrow \infty} P(|X_n - a| \geq \epsilon) = 0$$



Aire hachurée tend vers 0 quand  $n \rightarrow \infty$

## Loi faible des grands nombres

### Loi faible

- $(X_n)_{n \in \mathbb{N}}$  est une suite de variables aléatoires :
  - de même loi
  - d'espérance  $m$
  - possédant une variance  $\sigma^2$
  - **deux à deux** indépendantes
- alors la suite des variables  $\bar{X}_n = \frac{\sum_{k=1}^n X_k}{n}$  converge en probabilité vers  $m$

$\bar{X}_n$  est appelée **moyenne empirique**

$$E(\bar{X}_n) = m$$

$$V(\bar{X}_n) = \frac{\sigma^2}{n}$$

**conséquence** : échantillons de grandes tailles  $\implies$  bonne chance d'estimer  $m$

### Définition

- $(X_n)_{n \in \mathbb{N}}$  : suite de variables
- $a$  : constante
- $(X_n)$  **converge presque sûrement** vers  $a$  s'il y a une proba 1 que la suite des réalisations des  $X_n$  tende vers  $a$  :

$$P\left(\lim_{n \rightarrow \infty} X_n = a\right) = 1$$

 Définition plus exigeante que la convergence en probabilité

### Loi forte

- $(X_n)_{n \in \mathbb{N}}$  : suite de variables aléatoires
  - de même loi
  - d'espérance  $m$
  - possédant une variance  $\sigma^2$
  - **mutuellement** indépendantes
- alors la suite des variables  $\bar{X}_n = \frac{\sum_{k=1}^n X_k}{n}$  converge presque sûrement vers  $m$

**Interprétation** : échantillon de grande taille  
 $\implies$  bonne estimation de  $m$

### Définition

- $(X_n)_{n \in \mathbb{N}}$  : suite de variables
- $F_n$  : fonction de répartition de  $X_n$
- $X$  : variable de fonction de répartition  $F$
- La suite  $X_n$  **converge en loi** vers  $X$  lorsque  $F_n(x)$  tend vers  $F(x)$  en tout point de continuité de  $F$

**Notation** :  $X_n \xrightarrow{\text{loi}} X$

### Théorème central-limite

- $(X_n)_{n \in \mathbb{N}}$  : suite de variables
  - de même loi
  - d'espérance  $\mu$
  - de variance  $\sigma^2$
  - **mutuellement** indépendantes
- alors la suite des moyennes empiriques centrées réduites  $\frac{\bar{X}_n - \mu}{\sigma / \sqrt{n}}$  tend en loi vers la loi normale centrée réduite :

$$\frac{\bar{X}_n - \mu}{\sigma / \sqrt{n}} \xrightarrow{\text{loi}} \mathcal{N}(0, 1)$$

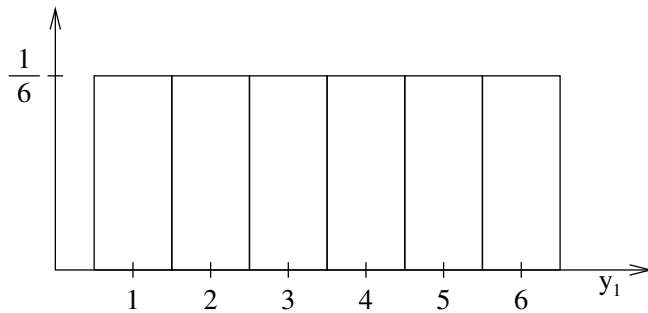
## Illustration du théorème central-limite (1/4)

Lancés de dés à 6 faces



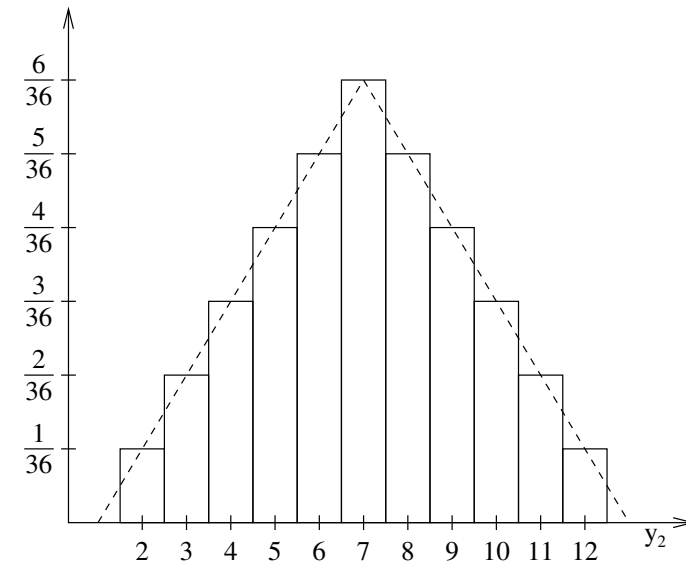
$$\Rightarrow \begin{cases} X_i = \text{résultat du jet du } i\text{ème dé} \\ \bar{X}_n = \text{somme des résultats des dés} \end{cases}$$

distribution de  $\bar{X}_n$  pour 1 jet de dé



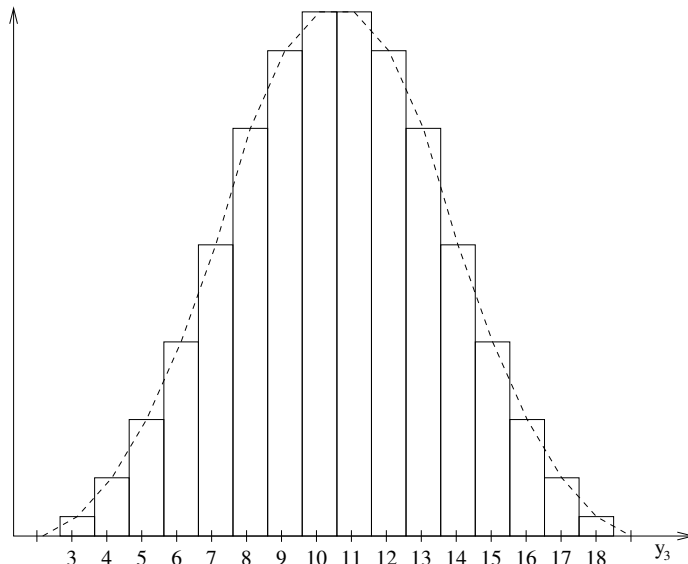
## Illustration du théorème central-limite (2/4)

distribution de  $\bar{X}_n$  pour 2 jets de dés



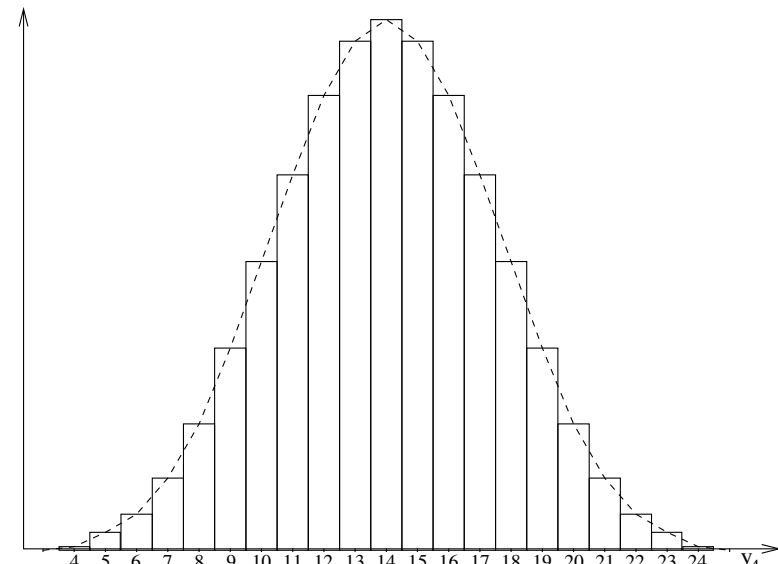
## Illustration du théorème central-limite (3/4)

distribution de  $\bar{X}_n$  pour 3 jets de dés



## Illustration du théorème central-limite (4/4)

distribution de  $\bar{X}_n$  pour 4 jets de dés



## Statistique inférentielle

### Hypothèses :

- les données = les observations
- observations  $(x_1, \dots, x_n)$  = réalisation d'une variable aléatoire multidimensionnelle  $X = (X_1, \dots, X_n)$
- observations = *échantillon empirique*
- l'échantillon est tiré suivant une loi  $X_0$
- chaque variable  $X_i$  suit la même loi que  $X_0$

### But :

- déduire des informations sur  $X_0$

### Exemples :

- sondages électoraux, études de marché
- tests de fiabilité/qualité
- bases de données  $\Rightarrow$  diagnostic

# Estimation ponctuelle

## Idée force

Dans moult études statistiques :

- population de grande taille  
 $\Rightarrow$  impossible de la connaître précisément
- possibilité de prélever des échantillons



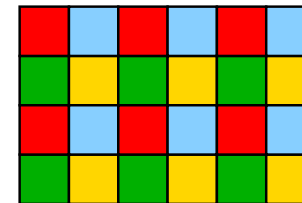
$\Rightarrow$  étudier l'échantillon et en déduire les caractéristiques de la population

**caractéristiques** : moyenne  $\mu$ , variance  $\sigma^2$ , proportion de succès  $p$

# Échantillonnage (1/2)

caractéristique de la population  $\Leftarrow$  échantillon

$\Rightarrow$  **problème** : comment prélever l'échantillon ?



## Définition d'un échantillon

**Échantillon** d'une population = sous-ensemble **représentatif** de la population

$\Rightarrow$  composition de l'échantillon due au hasard

### Prélèvement d'un échantillon

2 manières de prélever un échantillon aléatoirement :

- 1 **Échantillonnage avec remise** :  
choisir un individu au hasard  
noter la valeur de la variable d'intérêt pour celui-ci  
remettre l'individu dans la population  
réitérer le processus
- 2 **Échantillonnage sans remise** :  
même procédé mais sans remettre dans la population  
les individus sélectionnés

⚠ échantillonnage avec remise  $\implies$  un individu peut apparaître plusieurs fois dans l'échantillon

⚠ population de grande taille  $\implies$  avec remise  $\approx$  sans remise

⚠ résultats suivant valables uniquement pour des échantillons *avec remise*

$\implies$  on va travailler avec des échantillons i.i.d :

### Échantillon i.i.d

- échantillon de  $n$  individus
  - $X_i$  = variable aléatoire « valeur du  $i$ ème individu tiré »
  - les  $X_i$  sont mutuellement indépendants
  - les  $X_i$  ont tous la même distribution
- $\implies$  les  $X_i$  sont indépendants et identiquement distribués (i.i.d)

échantillons i.i.d  $\implies$  bonnes propriétés mathématiques

- $X$  variable aléatoire sur l'ensemble de la population
- espérance :  $\mu$ , variance :  $\sigma^2$
- échantillon de  $n$  individus  $\implies$  observation de  $n$  valeurs de  $X$
- $X_i$  : variable aléatoire correspondant au  $i$ ème individu
- échantillon i.i.d  $\implies$  espérance de  $X_i$  :  $\mu$ , variance de  $X_i$  :  $\sigma^2$

●  $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$  = moyenne de l'échantillon

**Problèmes** : que valent  $E(\bar{X})$  et  $V(\bar{X})$  ?

● 
$$E(\bar{X}) = \frac{1}{n} E(X_1 + X_2 + \dots + X_n)$$

$$= \frac{1}{n} (E(X_1) + E(X_2) + \dots + E(X_n))$$

$$= \frac{1}{n} (\mu + \mu + \dots + \mu) = \mu$$

variables  $X_i$  mutuellement indépendantes  $\implies$

● 
$$V(\bar{X}) = V\left(\frac{X_1 + X_2 + \dots + X_n}{n}\right) = \frac{1}{n^2} V(X_1 + X_2 + \dots + X_n)$$

$$= \frac{1}{n^2} (V(X_1) + V(X_2) + \dots + V(X_n))$$

$$= \frac{1}{n^2} (\sigma^2 + \sigma^2 + \dots + \sigma^2) = \frac{\sigma^2}{n}$$

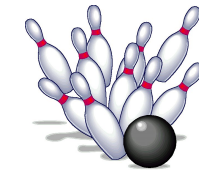
### Théorème 1

- $X$  : variable aléatoire
- espérance de  $X$  :  $\mu$ , variance de  $X$  :  $\sigma^2$
- échantillon de taille  $n$  avec remise sur  $X$
- $\bar{X}$  : variable aléatoire « moyenne de l'échantillon »
- Alors :  $E(\bar{X}) = \mu$  et  $V(\bar{X}) = \frac{\sigma^2}{n}$

### Corollaire

- $X$  : variable aléatoire suivant une loi normale  $\mathcal{N}(\mu; \sigma^2)$
- échantillon de taille  $n$  avec remise sur  $X$
- $\bar{X}$  : variable aléatoire « moyenne de l'échantillon »
- Alors :  $\bar{X} \sim \mathcal{N}(\mu; \frac{\sigma^2}{n})$

### Compétition en 2 phases :



- 1ère phase : 6 parties
- calcul du handicap :  
Handicap  $H = (215 - \text{moyenne des 6 parties}) \times 60\%$
- 2ème phase : 6 parties (score + le handicap  $H$ )

- Calcul du handicap  $\implies$  estimation de votre niveau
- 6 parties de la 1ère phase = échantillon de taille 6
- $X_i$  = variable aléatoire « score de la  $i$ ème partie »

idée : scores de l'échantillon  $\implies$  score moyen de la population

## Retour sur l'estimation de la moyenne

**Problème** : que faire si  $X$  ne suit pas une loi normale ?

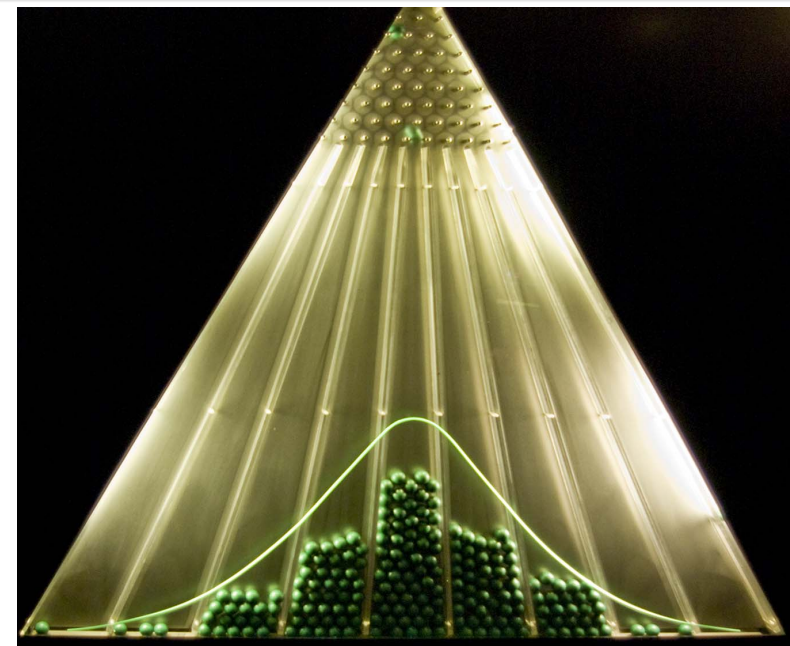
$\implies$  sauvé grâce au théorème central-limite :

### Théorème 2

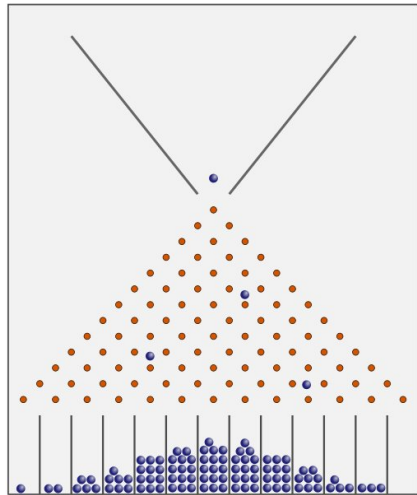
- $X$  : variable aléatoire
- espérance de  $X$  :  $\mu$ , variance de  $X$  :  $\sigma^2$
- échantillon de taille  $n$  avec remise sur  $X$
- $n$  suffisamment grand ( $n \geq 30$  si la distribution de  $X$  n'est pas trop dissymétrique,  $n \geq 50$  sinon)
- $\bar{X}$  : variable aléatoire « moyenne de l'échantillon »

- Alors :  $\frac{\bar{X} - E(\bar{X})}{\sqrt{V(\bar{X})}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim \mathcal{N}(0; 1)$

## Exemple 1 : la planche de Galton (1/2)



## Exemple 1 : la planche de Galton (2/2)



- chaque niveau  $\implies$  expérience de Bernoulli
- $\implies X \sim$  loi binomiale
- $\implies X \not\sim$  loi normale
- théorème 2  $\implies$

$$\frac{\bar{X} - E(\bar{X})}{\sqrt{V(\bar{X})}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim \mathcal{N}(0; 1)$$

## Exemple 2 : analyse des déchets

- Grenelle de l'environnement
  - $\implies$  réduction des déchets
  - $\implies$  analyse des déchets
- impossible à réaliser sur toute la population
  - $\implies$  échantillon de taille 100 :



450	320	320	390	410	415	380	390
440	350	400	380	430	400	375	...

- moyenne de l'échantillon = 390 kg/an/habitant
- $\bar{X}$  : variable aléatoire « moyenne de l'échantillon »
- $\sigma = 20$  supposé connu
- $\bar{X} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n}) \implies$  écart-type de  $\bar{X} = 20/10 = 2$

## Estimation de proportions de succès (1/3)

**Problème :**



- lancement de la Wii  $\implies$  étude de marché
- échantillon de taille  $n \implies$  succès = Wii, échec = PS3
- $p =$  proportion de succès dans toute la population
- $P_i$  : variable aléatoire « succès du  $i$ ème individu »

**Question :** peut-on déduire  $p$  en observant les  $p_i$  ?

## Estimation de proportions de succès (2/3)

**Question :** peut-on déduire  $p$  en observant les  $p_i$  ?

- $p =$  proportion de succès dans la population ( $p$  pas trop petit)
- échantillon i.i.d de taille  $n$  assez grand
- $P_i =$  variable aléatoire « succès du  $i$ ème individu »
- $P_i \sim$  loi binomiale  $\mathcal{B}(1, p)$
- $\bar{P} =$  moyenne de l'échantillon

Théorème central-limite  $\implies \frac{\bar{P} - E(\bar{P})}{\sqrt{V(\bar{P})}} \sim \mathcal{N}(0; 1)$

Or  $\bar{P} \sim \frac{1}{n}\mathcal{B}(n; p) \implies E(\bar{P}) = \frac{np}{n}$  et  $V(\bar{P}) = \frac{p(1-p)}{n}$



### Théorème 3

- échantillon i.i.d de taille  $n$  assez grand
- $p$  = proportion de succès dans la population ( $p$  pas trop petit)
- $\bar{P}$  = moyenne de l'échantillon

● Alors : 
$$\frac{\bar{P} - E(\bar{P})}{\sqrt{V(\bar{P})}} = \frac{\bar{P} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim \mathcal{N}(0; 1)$$

⇒ on peut estimer  $p$  en observant la valeur de  $\bar{P}$

Expert : 80%



v.s. 20%



- échantillon de 100 personnes ⇒ 70 Wii et 30 PS3

- Doit-on croire l'expert ?

si  $p = 80\%$  alors Théorème 3 ⇒ 
$$\frac{\bar{P} - 0,8}{\sqrt{\frac{0,8 \times 0,2}{100}}} = 25(\bar{P} - 0,8) \sim \mathcal{N}(0; 1)$$

$$\begin{aligned} \text{Prob}(\bar{P} \leq 0.7) &= \text{Prob}(25(\bar{P} - 0.8) \leq 25 \times (0.7 - 0.8)) \\ &= \text{Prob}(25(\bar{P} - 0.8) \leq -2,5) \approx 0,62\% \end{aligned}$$

## Application au réchauffement climatique

- Avril 2007 : étude tms-sofres / CNRS : opinion des gens sur le réchauffement climatique
- 1000 personnes de 15 ans et + interrogées ⇒ échantillon i.i.d
- 790 pensent qu'il y a un changement climatique
- 210 ne le pensent pas
- $\bar{P}$  : proportion de succès moyenne de l'échantillon
- $p$  : proportion de personnes pensant qu'il y a dérèglement climatique dans la population française



$$\frac{\bar{P} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim \mathcal{N}(0; 1)$$

⇒ estimation de  $p = 79\%$ , écart-type de  $\bar{P} \leq \sqrt{\frac{0,25}{1000}} \approx 0,015$

## Caractéristiques des estimateurs

Transparents précédents ⇒ si échantillon de grande taille alors  $\bar{X} \approx \mu \Rightarrow$  estimer  $\mu$  par  $\bar{X}$

### Estimateur non biaisé

- Estimateur  $T$  d'un paramètre  $\theta \Rightarrow$  valeur estimée  $\hat{\theta}$
- $T$  non biaisé si  $E(T) = \theta$



estimateur biaisé ⇒ on surévalue ou sous-évalue  $\theta$

### Estimateur convergent

- Estimateur  $T$  d'un paramètre  $\theta$
- $T$  convergent si  $E[(T - \theta)^2] \rightarrow 0$  lorsque la taille de l'échantillon ↗

moyenne  $\bar{X}$  et proportion de succès  $\bar{P}$  : estimateurs non biaisés et convergents



## Biais : estimation ponctuelle d'une variance (1/5)

- population = 4 nombres {1, 2, 3, 4}
- $\mu = \frac{5}{2}$  et  $\sigma^2 = \frac{5}{4}$
- échantillon de taille 2 avec remise :

échant.	espérance	variance	échant.	espérance	variance
1 1	$E(\bar{X}) = 1$	$V(\bar{X}) = 0$	1 2	$E(\bar{X}) = \frac{3}{2}$	$V(\bar{X}) = \frac{1}{4}$
1 3	$E(\bar{X}) = 2$	$V(\bar{X}) = 1$	1 4	$E(\bar{X}) = \frac{5}{2}$	$V(\bar{X}) = \frac{9}{4}$
2 1	$E(\bar{X}) = \frac{3}{2}$	$V(\bar{X}) = \frac{1}{4}$	2 2	$E(\bar{X}) = 2$	$V(\bar{X}) = 0$
2 3	$E(\bar{X}) = \frac{5}{2}$	$V(\bar{X}) = \frac{1}{4}$	2 4	$E(\bar{X}) = 3$	$V(\bar{X}) = 1$
3 1	$E(\bar{X}) = 2$	$V(\bar{X}) = 1$	3 2	$E(\bar{X}) = \frac{5}{2}$	$V(\bar{X}) = \frac{1}{4}$
3 3	$E(\bar{X}) = 3$	$V(\bar{X}) = 0$	3 4	$E(\bar{X}) = \frac{7}{2}$	$V(\bar{X}) = \frac{1}{4}$
4 1	$E(\bar{X}) = \frac{5}{2}$	$V(\bar{X}) = \frac{9}{4}$	4 2	$E(\bar{X}) = 3$	$V(\bar{X}) = 1$
4 3	$E(\bar{X}) = \frac{7}{2}$	$V(\bar{X}) = \frac{1}{4}$	4 4	$E(\bar{X}) = 4$	$V(\bar{X}) = 0$

- espérance  $\left\{ \begin{array}{l} \text{des moyennes des échantillons} = \mu \\ \text{des variances des échantillons} = \frac{5}{8} \neq \frac{5}{4} \implies \text{biais !} \end{array} \right.$

## Biais : estimation ponctuelle d'une variance (2/5)

- $X$  : variable aléatoire sur l'ensemble de la population
- $\sigma^2$  : variance de  $X$
- échantillon i.i.d. de taille  $n$
- $X_i$  : variable aléatoire correspondant au  $i$ ème individu
- $x_i$  : valeur observée de  $X_i$
- $\bar{X}$  : variable aléatoire « moyenne sur l'échantillon »
- $\bar{x}$  : valeur observée de  $\bar{X}$
- $s_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \left( \sum_{i=1}^n x_i^2 \right) - \bar{x}^2$
- $s_n^2$  = variance de l'échantillon

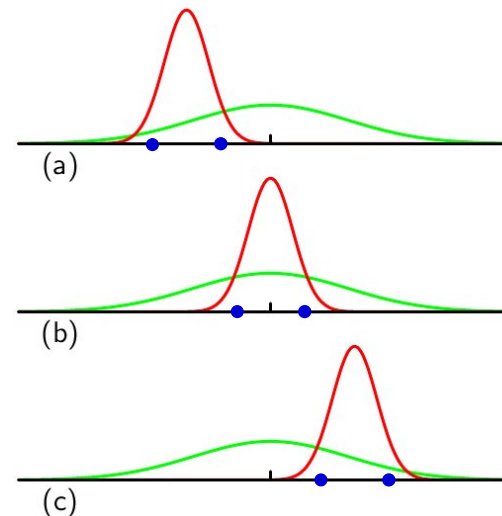
**Problème** : pourquoi  $s_n^2$  n'est-il pas un bon estimateur de  $\sigma^2$  ?

## Biais : estimation ponctuelle d'une variance (3/5)

- $S_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n} \left( \sum_{i=1}^n X_i^2 \right) - \bar{X}^2$
- $s_n^2$  est une réalisation de  $S_n^2$
- $E(S_n^2) = E \left( \frac{1}{n} \left( \sum_{i=1}^n X_i^2 \right) - \bar{X}^2 \right) = \frac{1}{n} \left( \sum_{i=1}^n E(X_i^2) \right) - E(\bar{X}^2)$
- Or échantillon i.i.d  $\implies \forall i, E(X_i^2) = E(X^2)$   
 $\implies E(S_n^2) = E(X^2) - E(\bar{X}^2)$
- Or  $V(X) = E(X^2) - E(X)^2$  et  $V(\bar{X}) = E(\bar{X}^2) - E(\bar{X})^2 = E(\bar{X}^2) - E(X)^2$   
 $\implies E(S_n^2) = V(X) + E(X)^2 - V(\bar{X}) - E(X)^2 = V(X) - V(\bar{X})$
- Or échantillon i.i.d  $\implies V(X) = \sigma^2$  et  $V(\bar{X}) = \frac{\sigma^2}{n}$

$$\implies E(S_n^2) = \sigma^2 - \frac{\sigma^2}{n} = \frac{n-1}{n} \sigma^2$$

## Biais : estimation ponctuelle d'une variance (4/5)



- courbe verte : la population
- courbes rouges : échantillons
- pts bleus : valeurs observées
- variance  $S_n^2$  sous-estimée : mesurée par rapport à la moyenne de l'échantillon au lieu de la moyenne de la population

### Variance corrigée

- $X$  : variable aléatoire sur la population
- $\sigma^2$  : variance de  $X$  sur cette population
- échantillon de taille  $n$  avec remise
- $X_i$  : variable aléatoire correspondant au  $i$ ème individu
- $\bar{X}$  : moyenne des  $X_i$
- alors :  $E \left( \frac{n}{n-1} \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{n} \right) = \sigma^2$
- **Variance corrigée** :  $\frac{n}{n-1}$  fois la variance de l'échantillon  $S_n^2$