

- Licence Sciences et technologie
- Mention science de la vie – L2



## INTRODUCTION A LA MODELISATION EN BIOLOGIE (BM1) - LU2SV382

### Modelisation statistique (Cours 1)

Martin LARSEN

[www.immulab.fr](http://www.immulab.fr)

- Licence Sciences et technologie
- Mention science de la vie – L2



## Prise en main d'Excel

Martin LARSEN

[www.immulab.fr](http://www.immulab.fr)

## Plan d'action

---

1. Statistiques descriptives
2. Visualisation des données
3. Distribution des données,
4. Intervalle de pari/confiance et Inférence statistique

## Excel - prise en main

---

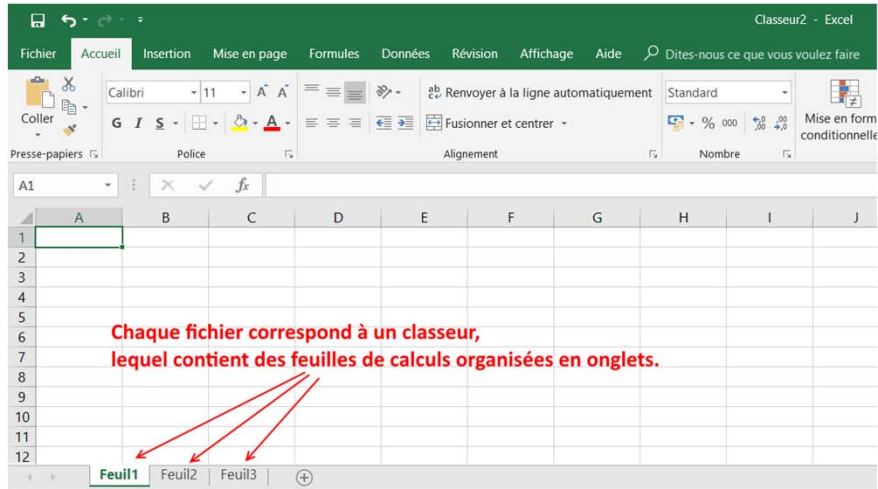
### Excel c'est quoi ?

Le logiciel Excel intègre des fonctions :

1. de calcul numérique (à l'aide de formules de calcul)
2. de représentation graphique
3. d'analyse de données (notamment de tableau croisé dynamique)
4. et de programmation, laquelle utilise les macros écrites dans un langage (Visual Basic) – **HORS CURSUS**

## Excel - prise en main

### Un fichier Excel c'est quoi ?



## Excel - prise en main

### Un tableau Excel c'est quoi ?

Colonnes : A, B, C ...G ...

Mois	CD musicaux	CD jeux	Livres	Magazines	Total par mois
Janvier	1 900	6 520	2 751	514	11 685
Février	1 888	4 522	2 987	687	10 085
Mars	1 457	4 523	1 237	125	7 342
Avril	1 255	9 563	1 588	385	12 790
Mai	957	7 521	2 430	873	11 781
Juin	1 250	6 385	2 100	268	10 004
Juillet	1 278	2 341	1 251	487	5 356
Août	958	4 588	1 955	453	7 953
Septembre	1 752	2 384	1 235	365	5 737
Octobre	1 521	5 367	2 840	569	10 298
Novembre	245	7 851	2 462	235	10 793
Décembre	1 238	4 853	3 100	751	9 943
<b>Total par produit</b>	<b>15 697</b>	<b>66 420</b>	<b>25 936</b>	<b>5 714</b>	<b>113 767</b>

Lignes : 1, 2, 3, 4, .....16, 17 ....

## Excel - prise en main

### Un tableau Excel c'est quoi ?

Zone de nom

La cellule B4 est la case située à l'intersection de la colonne B et de la ligne 4. Elle contient la valeur : 1900

2017: Vente d'articles en ligne					
Produit					
Mois	CD musicaux	CD jeux	Livres	Magazines	Total par mois
Janvier	1 900				1 900
Février					0
Mars					0
Avril					0
Mai					0
Juin					0
Juillet					0
Août					0
Septembre					0
Octobre					0
Novembre					0
Décembre					0
<b>Total par produit</b>	1 900	0	0	0	1 900

## Excel - prise en main

### Calcul numérique en Excel

	A	B	C	D
1	h/j	24		
2	min/h	60		
3	sec/min	60		
4				
5				
6				

## Excel - prise en main

### Calcul numérique en Excel

**Formule:**

Initié par =

Les opérateurs arithmétiques:

+ - \* / ^

Les opérateurs relationnels:

=, <>, <, >, <=, >=

Les opérateurs logiques:

ET, OU, NON (des fonctions en Excel)

	A	B	C	D
1	h/j	24		
2	min/h	60		
3	sec/min	60		
4				
5				
6				

Formula bar: =60\*60

## Excel - prise en main

### Calcul numérique en Excel

**Formule:**

Initié par =

Les opérateurs arithmétiques:

+ - \* / ^

	A	B	C	D
1	h/j	24		
2	min/h	60		
3	sec/min	60		
4	sec/h	3600		
5	sec/j	86400		
6				
7				
8				

Formula bar 1: =B1\*B2\*B3

Formula bar 2: =B2^2

Text: Ou : =B2\*B3

## Excel - prise en main

### Calcul numérique en Excel

#### Formule:

Initié par =

Les opérateurs arithmétiques:

+ - \* / ^

Les opérateurs relationnels:

=, <>, <, >, <=, >=

Les opérateurs logiques:

ET, OU, NON (des fonctions en Excel)

	A	B	C	D
1	h/j	24		
2	min/h	60		
3	sec/min	60		
4	sec/h	3600		
5	sec/j	86400		
6				
7				Plus grand : >
8				

	A	B	C	D	E
1	h/j	24			
2	min/h	60			
3	sec/min	60			
4	sec/h	3600			
5	sec/j	86400			
6					
7					ET

	A	B	C	D	E
1	h/j	24			
2	min/h	60			
3	sec/min	60			
4	sec/h	3600			
5	sec/j	86400			
6					
7					OU

	A	B	C	D	E
1	h/j	24			
2	min/h	60			
3	sec/min	60			
4	sec/h	3600			
5	sec/j	86400			
6					
7					NON

## Excel - prise en main

### Des fonctions en Excel

**Plage de cellules :** Reference 1 : Reference 2 (e.g. A1:A3)

Decaler(Reference 1, ligne, colonne) (e.g. A3 = Decaler(A1;2;0))

#### Statistique descriptive:

Moynne(Plage)

Mediane(Plage)

Mode(Plage)

Min(Plage)

Centile(Plage;0,25) (Q1)

Centile(Plage;0,75) (Q3)

Max(Plage)

Ecartype(Plage)=s

Var(Plage)=s<sup>2</sup> (echantillon)

#### Attention!

Var.P(Plage)=σ<sup>2</sup>(population)

#### Pour info!

Racine(4) = √4 = 2

$$s^2 = \frac{n}{n-1} \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n} = \frac{1}{n-1} \left( \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right)$$

$$\sigma^2 = \sum_{i=1}^n \frac{(x_i - \mu)^2}{n} = \frac{1}{n} \left( \sum_{i=1}^n x_i^2 - n\mu^2 \right)$$

## Excel - Diagramme en bâtons (superposés)

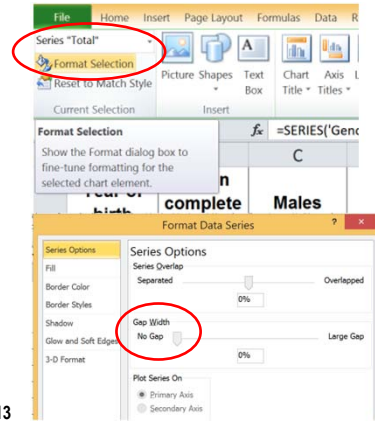
Démographie de la population française (Tableau 1):

	A	B	C	D	E
	Year of birth	Age in complete d years	Males	Females	Total
1					
2	2018	0	364,155	347,749	711,904
3	2017	1	370,453	355,472	725,925
4	2016	2	378,518	363,162	741,680
5	2015	3	387,906	372,402	760,308
6	2014	4	399,232	387,042	786,274
7	2013	5	407,611	389,920	797,531
8	2012	6	417,471	396,835	814,306
9	2011	7	418,623	403,349	821,972
10	2010	8	429,919	412,555	842,474
11	2009	9	427,917	408,232	836,149
92	1928	90	55,382	124,322	179,704
93	1927	91	44,797	105,456	150,253
94	1926	92	34,519	91,072	125,591
95	1925	93	27,317	76,447	103,764
96	1924	94	20,525	61,235	81,760
97	1923	95	14,477	48,398	62,875
98	1922	96	10,101	37,882	47,983
99	1921	97	7,239	27,754	34,993
100	1920	98	4,977	19,813	24,790
101	1919	99	2,058	8,273	10,331
102	1918 or before	100 or over	2,976	12,670	15,646
103		<b>Total</b>	<b>32,394,531</b>	<b>34,598,168</b>	<b>66,992,699</b>

Source: <https://www.insee.fr/en/statistiques/2382597?sommaire=2382613>



1. Affichez le nombre d'individus pour chaque groupe d'âge sous forme de diagramme en bâtons.
2. Formater les données - supprimer les écart entre des bâtons.



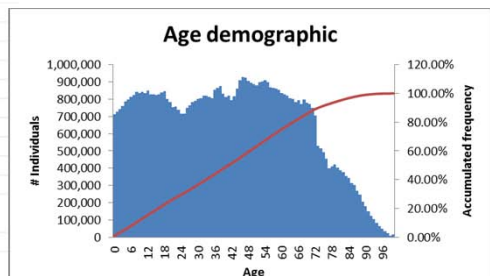
## Excel - Diagramme en bâtons superposés

Démographie de la population française (Tableau 1):

	A	B	C	D	E
	Year of birth	Age in complete d years	Males	Females	Total
1					
2	2018	0	364,155	347,749	711,904
3	2017	1	370,453	355,472	725,925
4	2016	2	378,518	363,162	741,680
5	2015	3	387,906	372,402	760,308
6	2014	4	399,232	387,042	786,274
7	2013	5	407,611	389,920	797,531
8	2012	6	417,471	396,835	814,306
9	2011	7	418,623	403,349	821,972
10	2010	8	429,919	412,555	842,474
11	2009	9	427,917	408,232	836,149
92	1928	90	55,382	124,322	179,704
93	1927	91	44,797	105,456	150,253
94	1926	92	34,519	91,072	125,591
95	1925	93	27,317	76,447	103,764
96	1924	94	20,525	61,235	81,760
97	1923	95	14,477	48,398	62,875
98	1922	96	10,101	37,882	47,983
99	1921	97	7,239	27,754	34,993
100	1920	98	4,977	19,813	24,790
101	1919	99	2,058	8,273	10,331
102	1918 or before	100 or over	2,976	12,670	15,646
103		<b>Total</b>	<b>32,394,531</b>	<b>34,598,168</b>	<b>66,992,699</b>

Source: <https://www.insee.fr/en/statistiques/2382597?sommaire=2382613>

1. Affichez le nombre d'individus pour chaque groupe d'âge sous forme de diagramme en bâtons.
2. Formater les données - supprimer les écart entre des bâtons.
3. Calcule la fréquence accumulée et trace sous forme de ligne. (en cellule F2: =Somme(\$E\$2:E2) et puis copie).  
!! Touche de fonction F4 transforme E2 en \$E\$2.  
\$ Fixe le référence – colonne, ligne ou les deux.



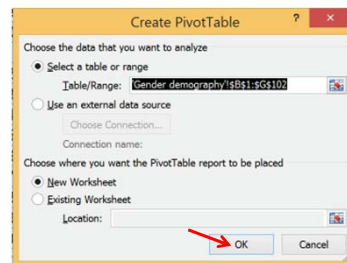
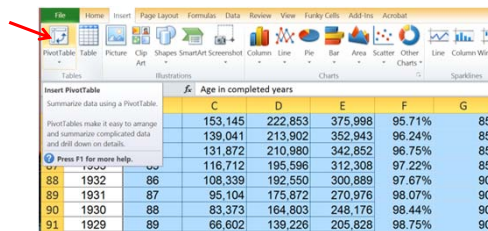
## Excel - Tableau Croisé Dynamique (Pivot tables)

Démographie de la population française (Tableau 1):

Year of birth	Age in completed years	Males	Females	Total	Freq Acc	<Age
2018	0	364,155	347,749	711,904	1.06%	5
2017	1	370,453	355,472	725,925	2.15%	5
2016	2	376,518	363,162	741,680	3.25%	5
2015	3	387,906	372,402	760,308	4.39%	5
2014	4	399,232	387,042	786,274	5.58%	5
2013	5	407,611	389,920	797,531	6.75%	10
1922	96	10,101	37,882	47,983	99.87%	100
1921	97	7,239	27,754	34,993	99.92%	100
1920	98	4,977	19,813	24,790	99.96%	100
1919	99	2,058	8,273	10,331	99.98%	100
1918 or before	100+	2,976	12,670	15,646	100.06%	#VALUE!

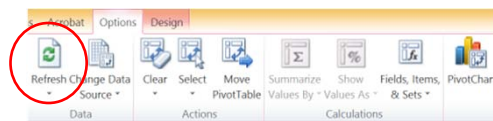
1. Faire des groupes d'âge avec un intervalle de n années:  

$$= n * \text{Plafond} ((B2+1) / n; 1)$$
 Par exemple. n = 5
2. Corrigez les lignes 102 à 100.
3. Sélectionner des données et créer un tableau croisé dynamique

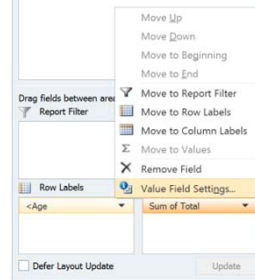
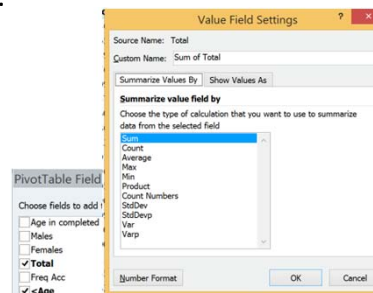


## Excel - Tableau Croisé Dynamique (Pivot tables)

Démographie de la population française (Tableau 1):



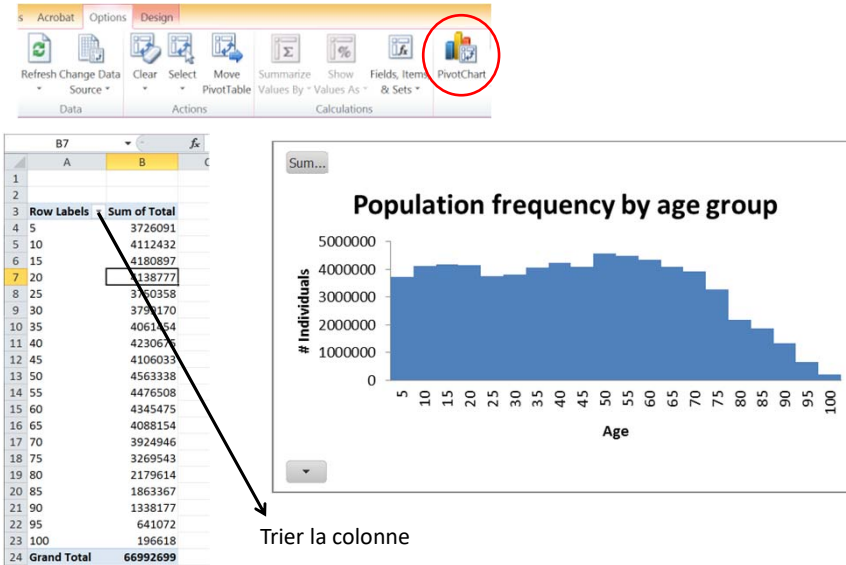
Row Labels	Sum of Total
5	3726091
10	4112432
15	4180897
20	4138777
25	3750358
30	3799170
35	4061454
40	4230675
45	4106033
50	4563338
55	4476508
60	4345475
65	4088154
70	3924946
75	3269543
80	2179614
85	1863367
90	1338177
95	641072
100	196618
<b>Grand Total</b>	<b>66992699</b>





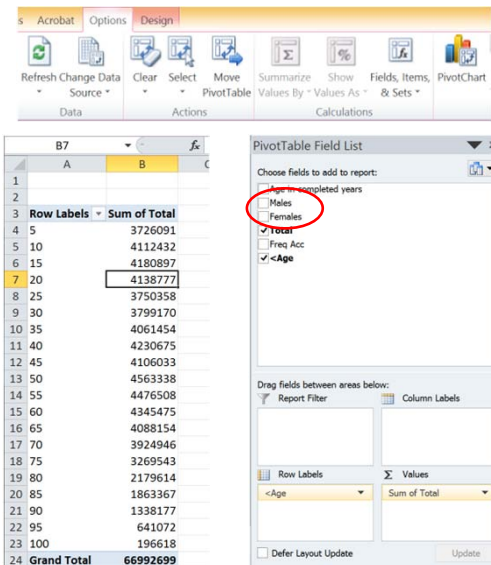
## Excel - Tableau Croisé Dynamique (Pivot tables)

Démographie de la population française (Tableau 1):



## Excel - Tableau Croisé Dynamique (Pivot tables)

Démographie de la population française (Tableau 1):

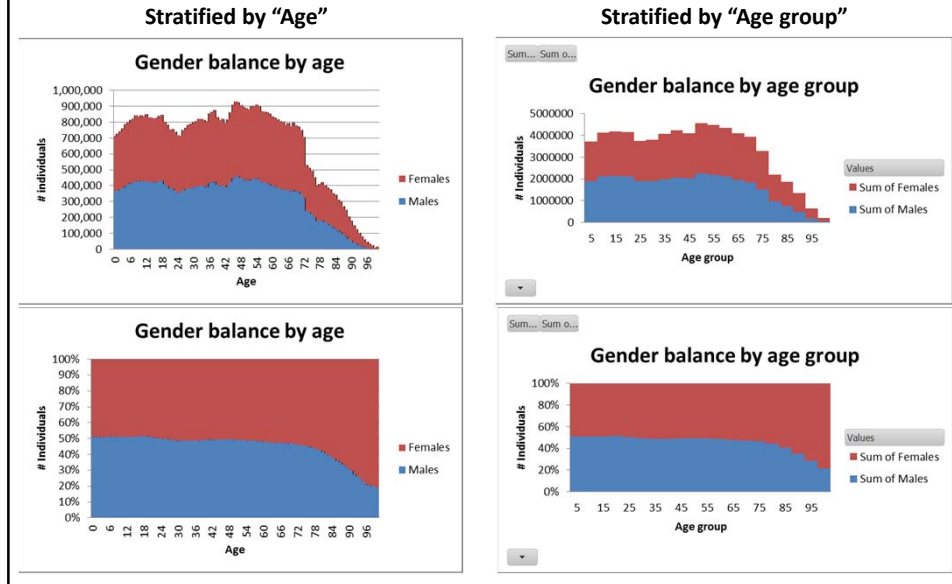


Analyser la parité hommes-femmes stratifiée par « âge » et par « groupes d'âge ».

1. Utiliser le jeu de données complet (« age »)
2. Utilisez un tableau croisé dynamique pour effectuer l'analyse stratifiée par groupes d'âge.

## Excel - Tableau Croisé Dynamique (Pivot tables)

Démographie de la population française (Tableau 1):



## Population / échantillon



Population

Loi connue dans la population

- Loi Binomiale
- Loi Normale
- ...



Échantillon

Prédiction de ce qui sera observé dans l'échantillon

⇒ Intervalle de pari

## Population / échantillon



Population



Généralisation / Inférence



Échantillon



Observation



Estimation





# Moyenne et pourcentage



← Intervalle de confiance

→ Intervalle de Pari

## Intervalle de Pari

- $IP_{1-\alpha} = \left[ \mu - u_\alpha \frac{\sigma}{\sqrt{n}}; \mu + u_\alpha \frac{\sigma}{\sqrt{n}} \right]$
- $IP_{1-\alpha} = \left[ \pi - u_\alpha \sqrt{\frac{\pi(1-\pi)}{n}}; \pi + u_\alpha \sqrt{\frac{\pi(1-\pi)}{n}} \right]$

• A partir de la connaissance de la population on imagine ce que l'on obtiendra dans un échantillon

## Intervalle de confiance

- $IC_{1-\alpha} = \left[ m - u_\alpha \frac{s}{\sqrt{n}}; m + u_\alpha \frac{s}{\sqrt{n}} \right]$
- $IC_{1-\alpha} = \left[ p - u_\alpha \sqrt{\frac{p(1-p)}{n}}; p + u_\alpha \sqrt{\frac{p(1-p)}{n}} \right]$

• A partir de l'observation d'un échantillon on en déduit une information sur la population



Deux échantillons viennent-ils d'une ou de deux populations ?





## Démarche



- Construire un intervalle de confiance à partir des données observées et voir si cela correspond à ce que l'on attend en population

Exemple : Le traitement A est-il plus efficace que le traitement B ?

- Estimer l'efficacité avec les deux traitements
- Calculer la différence et l'intervalle de confiance de cette différence
- Hypothèse : Si  $A = B$  alors la différence en population  $\Delta = 0$
- Voir si l'intervalle de confiance avec un certain risque est compatible avec cette hypothèse



## Démarche



- Construire un intervalle de pari à partir des hypothèses et voir si les données de l'échantillon correspondent à ce qui est attendu

Exemple : Le traitement A est-il plus efficace que le traitement B ?

- Hypothèses :
  - Hypothèse nulle  $H_0$  : Si  $A = B$  alors la différence en population  $\Delta = 0$
  - Hypothèse alternative  $H_1$  : Si  $A \neq B$  alors la différence en population  $\Delta \neq 0$
- Trouver une statistique de test correspondant à cette hypothèse
- Etablir un intervalle de pari de cette statistique
- Estimer l'efficacité avec les deux traitements
- Voir si l'intervalle de pari est compatible avec l'estimation

## Distribution de X et $\bar{X}$

X est une variable aléatoire continue (VAC)

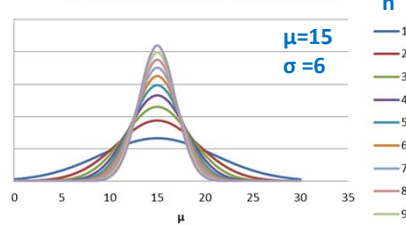
➤ Distribution de Echantonnage de Moyenne:

$$X \sim \mathcal{N}(\mu_x, \sigma_x) \Rightarrow \bar{X} \sim \mathcal{N}\left(\mu_x, \frac{\sigma_x}{\sqrt{n}}\right)$$

Pour  $n=1$ , nous avons  $X \sim \mathcal{N}(\mu_x, \sigma_x) = \bar{X} \sim \mathcal{N}\left(\mu_x, \frac{\sigma_x}{\sqrt{n}}\right)$

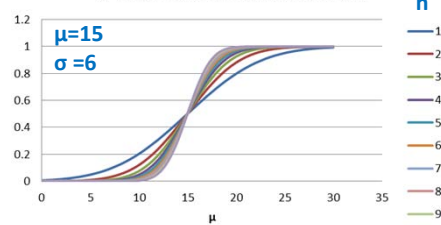
Fonction de densité de probabilité

$\bar{X}$  distribution for different n



Probabilité accumuler

$\bar{X}$  distribution for different n



## Distribution de X et $\bar{X}$

LOI.NORMALE.N(x;μ;σ/RACINE(N);cumulative)

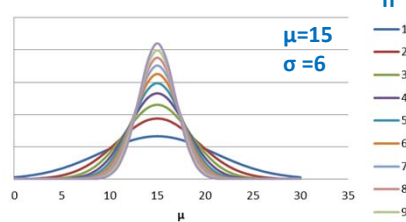
	A	B	C	D	E	F	G	H
mu		15						
sigma		6						
n		1	2	3	4	5	6	7
X	1	2	3	4	5	6	7	
	0	0.002921383	0.000182	9.77E-06	4.96E-07	2.43E-08	1.17E-09	5.56E-11
	1	0.004370315	0.000406	3.27E-05	2.48E-06	1.82E-07	1.31E-08	9.32E-10
	2	0.006358771	0.00086	0.000101	1.11E-05	1.19E-06	1.25E-07	1.29E-08
	3	0.008998494	0.001722	0.000285	4.46E-05	6.75E-06	1E-06	1.46E-07
	4	0.012385194	0.003263	0.000744	0.00016	3.33E-05	6.8E-06	1.37E-06
	5	0.016579523	0.005847	0.001785	0.000514	0.000143	3.91E-05	1.05E-05
	6	0.021586266	0.009911	0.003941	0.001477	0.000536	0.000191	6.69E-05
	7	0.027335012	0.015893	0.008002	0.003799	0.001746	0.000786	0.000349

Pour  $n=1$ , nous avons  $X \sim \mathcal{N}(\mu_x, \sigma_x) = \bar{X} \sim \mathcal{N}\left(\mu_x, \frac{\sigma_x}{\sqrt{n}}\right)$

Fonction de densité de probabilité

Cumulative = FAUX

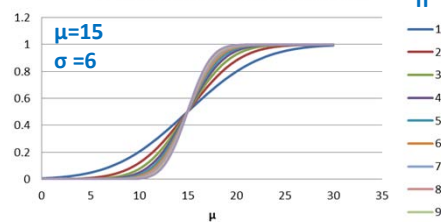
$\bar{X}$  distribution for different n



Probabilité accumuler

Cumulative = VRAI

$\bar{X}$  distribution for different n



## Distribution de X et $\bar{X}$

X est une variable aléatoire continue (VAC)

➤ Distribution de Echantonnage de Moyenne:

$$X \sim \mathcal{N}(\mu_x, \sigma_x) \Rightarrow \bar{X} \sim \mathcal{N}\left(\mu_x, \frac{\sigma_x}{\sqrt{n}}\right)$$

➤ Théorème centrale limite (TCL):

X suit n'importe quelle distribution. Pour  $n > 30 \Rightarrow \text{TCL} \Rightarrow \bar{X} \sim \mathcal{N}\left(\mu_x, \frac{\sigma_x}{\sqrt{n}}\right)$

➤ Autre

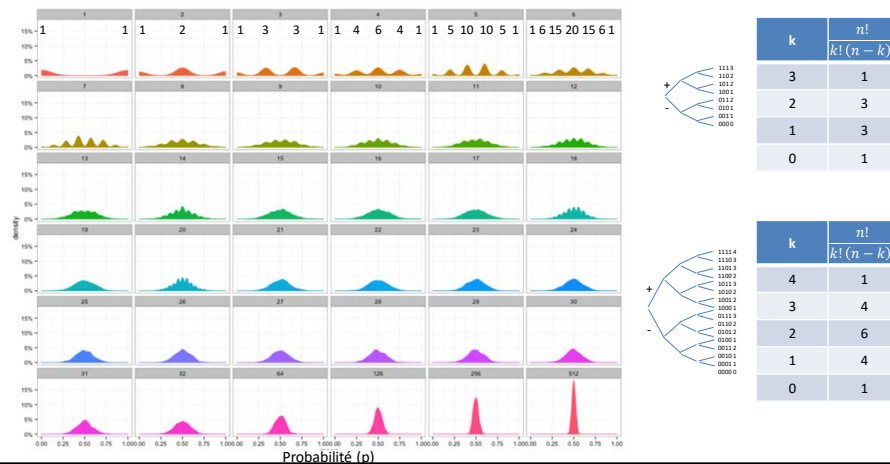
En l'absence de distribution connue pour X et  $n < 30$ , nous ne pouvons pas identifier une distribution pour  $\bar{X} \sim ??$

## Théorème centrale limite

**E.g. distribution binomial :**  $X \sim B(n, \pi) \Rightarrow \left[ \begin{array}{l} n\pi > 5 \\ n(1-\pi) > 5 \end{array} \right] \Rightarrow N(\mu = n\pi, \sigma^2 = n\pi(1-\pi))$   
 $\pi = 0.5$  and  $n \in [1-512]$   $P(X = k) = \binom{n}{k} p^k (1-p)^{n-k} = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$

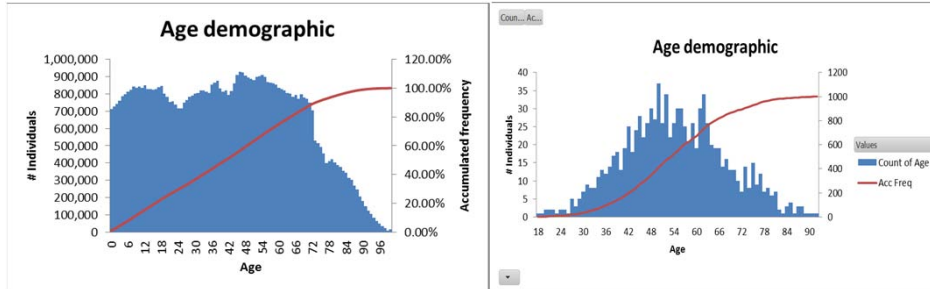
Fréquence de succès au cours de n épreuves de Bernoulli

Le moyenne de 1000 échantillons était analyse



## Distribution d'Age (statistique descriptive)

Démographie de la population française (Tableau 1) et une échantillon inconnu (Tableau 2):



Descriptive statistics of population:

n	66,992,699
$\mu$	41.2 years
$\sigma^2$	586.4 years <sup>2</sup>
$\sigma$	24.2 years

Descriptive statistics of sample

n	1000	=NB(A2:A1001)
Mean	54.2 years	=moyenne(A2:A1001)
$s^2$	184.3 years <sup>2</sup>	=Var.S(A2:A1001)
s	13.6 years	=ECARTYPE(A2:A1001)



## Moyenne et pourcentage



← Intervalle de confiance

Intervalle de Pari →

### Intervalle de Pari

- $IP_{1-\alpha} = \left[ \mu - u_{\alpha} \frac{\sigma}{\sqrt{n}}; \mu + u_{\alpha} \frac{\sigma}{\sqrt{n}} \right]$
- $IP_{1-\alpha} = \left[ \pi - u_{\alpha} \sqrt{\frac{\pi(1-\pi)}{n}}; \pi + u_{\alpha} \sqrt{\frac{\pi(1-\pi)}{n}} \right]$
- A partir de la connaissance de la population on imagine ce que l'on obtiendra dans un échantillon

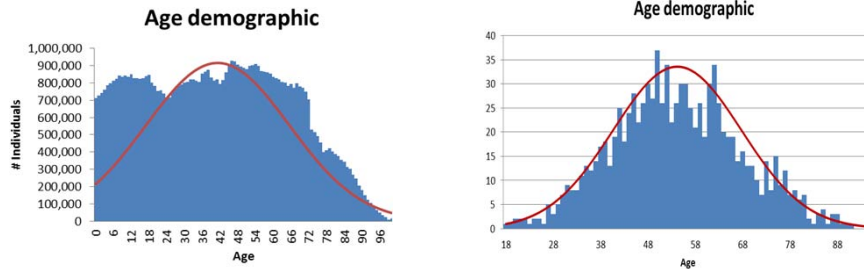
### Intervalle de confiance

- $IC_{1-\alpha} = \left[ m - u_{\alpha} \frac{s}{\sqrt{n}}; m + u_{\alpha} \frac{s}{\sqrt{n}} \right]$
- $IC_{1-\alpha} = \left[ p - u_{\alpha} \sqrt{\frac{p(1-p)}{n}}; p + u_{\alpha} \sqrt{\frac{p(1-p)}{n}} \right]$
- A partir de l'observation d'un échantillon on en déduit une information sur la population



## Intervalle de pari vs. Intervalle de confiance

Démographie de la population française (Tableau 1) et une échantillon inconnu (Tableau 2):



Statistique descriptive (population):

n 66,992,699  
 $\mu$  41.2 years  
 $\sigma^2$  586.4 years<sup>2</sup>  
 $\sigma$  24.2 years

X ne suis pas  $\mathcal{N}(\mu, \sigma)$  (Visuellement)  
 Mais Théorème Central Limite (TCL,  $n > 30$ ):

$$\bar{X} \sim \mathcal{N}\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

$$IP_{1-\alpha}: \bar{x} \in \left[\mu - u_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \mu + u_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right] \leftarrow \text{Intervalle rétréci}$$

$N \rightarrow \infty$   $\nearrow$

Statistique descriptive (échantillon)

n 1000  
 Mean 54.2 years  
 $s^2$  184.3 years<sup>2</sup>  
 s 13.6 years

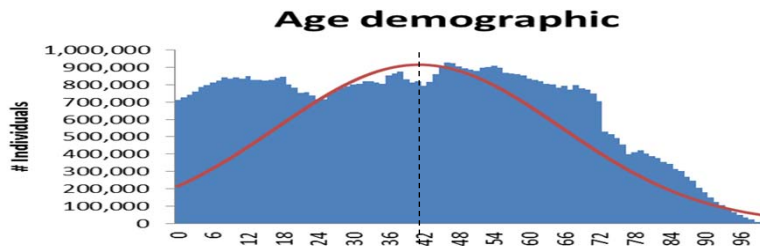
$$IC_{1-\alpha}: \mu \in \left[\bar{x} - u_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + u_{\alpha/2} \frac{s}{\sqrt{n}}\right]$$

$\nwarrow$

= [53.3; 55.0] avec risque  $\alpha=5\%$   
 Valider précision de s:  
 précise ( $n > 30$ ) (loi normale)  
 Peu- précise ( $n < 30$ ) (loi de Student)

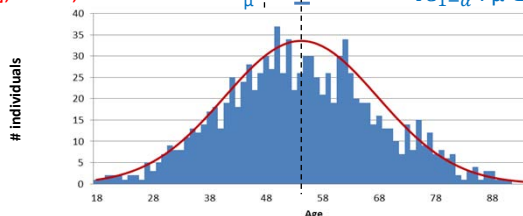
## Intervalle de pari vs. Intervalle de confiance

Démographie de la population française (Tableau 1) et une échantillon inconnu (Tableau 2):



$$IP_{1-\alpha}: \bar{x} \in [39.7; 42.7], \alpha=5\%, n=1000$$

$$IC_{1-\alpha}: \mu \in [53.3; 55.0], \alpha=5\%, n=1000$$



**Conclusion:**

- $\mu \notin IC_{1-\alpha}$ ; **Inférence:** l'échantillon a été sélectionné au hasard dans une autre population avec  $\mu \in [53.3; 55.0]$  au risque 5%
- $\bar{x} \notin IP_{1-\alpha}$ ; **Prédiction:** l'échantillon n'a pas été sélectionné au hasard dans la population avec  $\mu = 41.2$  au risque 5%

## Intervalle de pari vs. Intervalle de confiance

Démographie de la population française (Tableau 1) et une échantillon inconnu (Tableau 2):

