

# TP3: Intervalle de confiance

Tanguy ROUDAUT — Tadios QUINIO

FIPASE 24

22 Septembre 2022

## 1 Dissection d'une gaussienne

Données mesurées : 0.82 0.87 0.77 0.96 0.75 0.83 0.87 0.81

**Question 1 :** Considérant que l'échantillon a été engendré par une loi gaussienne, donner un intervalle de confiance pour son espérance. On utilisera les fonctions `tinv`, `mean` et `var` ou `std/t.ppf`, `np.mean`, `np.std`. Préciser toutes les hypothèses que vous reprenez.

Dans cette question, l'échantillon est de longueur  $n = 8$ .

Si un échantillon est de longueur  $n < 30$  avec une  $\sigma^2$  inconnu, il faut utiliser le cas suivant : *Intervalles de confiance de la moyenne d'une distribution normale,  $\sigma^2$  inconnue*

Ce qui nous amène à utiliser les formules 1 et 2 pour déterminer l'intervalle de confiance  $[I, u]$ .

$$I = \overline{X}_n - t_{n-1, \frac{\alpha}{2}} \frac{S_{n-1}}{\sqrt{n}} \quad (1) \quad \left| \quad u = \overline{X}_n + t_{n-1, \frac{\alpha}{2}} \frac{S_{n-1}}{\sqrt{n}} \quad (2)$$

Comme  $\sigma^2$  est inconnue, on peut l'approcher par  $S_{n-1}^2$  puis en déduire  $S_{n-1}$ .

On peut également rappeler que  $t_{n-1, \frac{\alpha}{2}}$  est le quantile d'ordre  $\frac{\alpha}{2}$  d'une loi de Student de  $n - 1$  degrés de liberté, déterminable en python avec `scipy.stats.t.ppf()`.

```
1 x = np.array([0.82, 0.87, 0.77, 0.96, 0.75, 0.83, 0.87, 0.81])
2 xn = np.mean(x)
3 ecart_type_x = np.std(x, ddof=1)
4 n = len(x)
5 ic = 95
6 alpha = 1 - ic / 100
7 t = stats.t.ppf(1 - (alpha / 2), n - 1)
8
9 I = xn - (t * ecart_type_x) / np.sqrt(n)
10 u = xn + (t * ecart_type_x) / np.sqrt(n)
11
12 I = round(I, 3)
13 u = round(u, 3)
14
15 print("Borne inférieur:", I, "\n", "Borne supérieur:", u)
```

Listing 1 – Code Python question 1

```
1 Borne inférieur: 0.779
2 Borne supérieur: 0.890
```

Listing 2 – Résultat du code

On peut donc conclure qu'il y'a 95% de chance que l'espérance se trouve dans l'intervalle  $[0,779; 0,890]$ .

**Question 2 :** Les données sont maintenant [0.84 0.87 0.89 0.73 0.84 0.81 0.88 0.85 0.89 0.79 0.79 0.90 0.59 0.75 0.67 0.76 0.86 0.88 0.70 0.75 0.81 0.77 0.83 0.84 0.71 0.78 0.59 0.91 0.74 0.68 0.77 0.66 0.80 0.74 1.02 0.91 0.55 0.84 0.66 0.77]. Considérant que l'échantillon a été engendré par une loi gaussienne, donner un intervalle de confiance pour son espérance. On utilisera les fonctions (norminv, mean et var ou std ou norm.ppf, np.mean, np.std). Préciser toutes les hypothèses que vous reprenez. Quel est dans ce cas, l'intérêt de la loi gaussienne ?

Dans cette question l'échantillon est de longueur  $n = 40$ .

Pour un échantillon de longueur  $n > 30$  suivant une distribution normale centrée réduite  $\mathcal{N}(0, 1)$  avec  $\sigma^2$  inconnu, il faut utiliser le cas suivant : *Intervalle de confiance de la moyenne d'une distribution dans le cas d'un grand échantillon*

Ce qui nous amène à utiliser la formule 3 pour déterminer l'intervalle de confiance.

$$\hat{p} - z_{\frac{\alpha}{2}} \sqrt{\frac{S_{n-1}}{n}} \leq p \leq \hat{p} + z_{\frac{\alpha}{2}} \sqrt{\frac{S_{n-1}}{n}} \quad (3)$$

$$\text{avec} \quad z_{\frac{\alpha}{2}} = \text{erf}^{-1}\left(1 - \frac{\alpha}{2}\right) \quad \text{et} \quad \alpha = 1 - \frac{IC\%}{100}$$

Comme  $\sigma^2$  est inconnue, on peut l'approcher par  $S_{n-1}^2$  puis en déduire  $S_{n-1}$ .

On peut également rappeler que  $z_{\frac{\alpha}{2}}$  peut être déterminé avec la table et en python avec `stats.norm.ppf()`.

```
1 x = np.array([0.84, 0.87, 0.89, 0.73, 0.84, 0.81, 0.88, 0.85, 0.89, 0.79, 0.79, 0.90,
2               0.59, 0.75, 0.67, 0.76, 0.86, 0.88, 0.70, 0.75, 0.81, 0.77, 0.83, 0.84,
3               0.71, 0.78, 0.59, 0.91, 0.74, 0.68, 0.77, 0.66, 0.80, 0.74, 1.02, 0.91,
4               0.55, 0.84, 0.66, 0.77])
5
6 p = np.mean(x)
7 ecart_type_x = np.std(x, ddof=1)
8 n = len(x)
9 ic = 95
10 alpha = 1 - ic / 100
11 z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
12
13 borne_inf = p - z * (ecart_type_x/np.sqrt(n))
14 borne_sup = p + z * (ecart_type_x/np.sqrt(n))
15
16 borne_inf = round(borne_inf, 3)
17 borne_sup = round(borne_sup, 3)
18
19 print("Borne inférieure:", borne_inf, "\n", "Borne supérieure:", borne_sup)
```

Listing 3 – Code Python question 2

```
1 Borne inférieure: 0.755
2 Borne supérieure: 0.816
```

Listing 4 – Résultat du code

On peut donc conclure qu'il y'a 95% de chance que l'espérance se trouve dans l'intervalle [0,755;0,816]

## 2 Recoder l'exercice 4 du TD : sondages

À la veille d'une consultation électorale, nous effectuons un sondage.

**Question 3 :** Dans un échantillon représentatif de 1000 personnes, 500 personnes déclarent vouloir voter pour Dupond, 250 pour Durand et 50 pour Duroc. Donner les intervalles de confiance à 95% et 99% du pourcentage de personnes ayant l'intention de voter Dupond, Durand ou Duroc.

Formules utilisées :

$$\hat{p} - z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \quad (4)$$

$$\text{avec} \quad z_{\frac{\alpha}{2}} = \text{erf}^{-1}\left(1 - \frac{\alpha}{2}\right) \quad \text{et} \quad \alpha = 1 - \frac{IC\%}{100}$$

```

1 n = 1000
2 n_dupond = 500
3 n_durand = 250
4 n_duroc = 50
5
6 def intervalle_confiance(n, x, ic ):
7     p = x/n
8     alpha = 1 - ic / 100
9     z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
10    borne_inf = p - z * (np.sqrt(p*(1-p)/n))
11    borne_sup = p + z * (np.sqrt(p*(1-p)/n))
12
13    return round(borne_inf, 3), round(borne_sup, 3)
14
15 borne_inf_dupond_95, borne_sup_dupond_95 = intervalle_confiance(n, n_dupond, 95)
16 borne_inf_durand_95, borne_sup_durand_95 = intervalle_confiance(n, n_durand, 95)
17 borne_inf_duroc_95, borne_sup_duroc_95 = intervalle_confiance(n, n_duroc, 95)
18 borne_inf_dupond_99, borne_sup_dupond_99 = intervalle_confiance(n, n_dupond, 99)
19 borne_inf_durand_99, borne_sup_durand_99 = intervalle_confiance(n, n_durand, 99)
20 borne_inf_duroc_99, borne_sup_duroc_99 = intervalle_confiance(n, n_duroc, 99)
21
22 print(" Dupond")
23 print("\t --> Borne inférieur à 95%:", borne_inf_dupond_95)
24 print("\t --> Borne supérieur à 95%:", borne_sup_dupond_95)
25 print("\t\t\t -----")
26 print("\t --> Borne inférieur à 99%:", borne_inf_dupond_99)
27 print("\t --> Borne supérieur à 99%:", borne_sup_dupond_99)
28 print(" Durand")
29 print("\t --> Borne inférieur à 95%:", borne_inf_durand_95)
30 print("\t --> Borne supérieur à 95%:", borne_sup_durand_95)
31 print("\t\t\t -----")
32 print("\t --> Borne inférieur à 99%:", borne_inf_durand_99)
33 print("\t --> Borne supérieur à 99%:", borne_sup_durand_99)
34 print(" Duroc")
35 print("\t --> Borne inférieur à 95%:", borne_inf_duroc_95)
36 print("\t --> Borne supérieur à 95%:", borne_sup_duroc_95)
37 print("\t\t\t -----")
38 print("\t --> Borne inférieur à 99%:", borne_inf_duroc_99)
39 print("\t --> Borne supérieur à 99%:", borne_sup_duroc_99, end="\n\n")

```

Listing 5 – Code Python question 3

```

1 Dupond
2   --> Borne inférieur à 95%: 0.469
3   --> Borne supérieur à 95%: 0.531
4   -----
5   --> Borne inférieur à 99%: 0.459
6   --> Borne supérieur à 99%: 0.541

```

```

7 Durand
8 --> Borne inférieur à 95%: 0.223
9 --> Borne supérieur à 95%: 0.277
10 -----
11 --> Borne inférieur à 99%: 0.215
12 --> Borne supérieur à 99%: 0.285
13 Duroc
14 --> Borne inférieur à 95%: 0.036
15 --> Borne supérieur à 95%: 0.064
16 -----
17 --> Borne inférieur à 99%: 0.032
18 --> Borne supérieur à 99%: 0.068

```

Listing 6 – Résultat du code

**Question 4 :** Nous évaluons le pourcentage de personnes ayant l'intention de voter pour un quatrième candidat, Duval, à 17%. Combien faut-il interroger de personnes pour obtenir une précision de 1% pour l'intervalle de confiance (à 95%) de la proportion de personnes ayant l'intention de voter Duval ?

$$n = \left( \sqrt{\frac{Z_{\alpha/2}}{E}} \right)^2 \hat{p}(1 - \hat{p}) \quad \text{avec } E = \text{erreur} = \text{précision} \quad (5)$$

```

1 ic = 95
2 alpha = 1 - ic / 100
3 p = 17/100
4 err = 1/100
5 z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
6 n = (z/err)**2 * p*(1-p)
7 n = np.ceil(n)
8
9 print("Nombre de personne à interroger:", n, end='\n\n')

```

Listing 7 – Code Python question 4

```

1 Nombre de personne à interroger: 5421.0

```

Listing 8 – Résultat du code

### 3 Dommages de casques

Un échantillon aléatoire de 50 casques de motos et de courses automobiles a été testé à l'impact. Des dommages ont été observés pour 18 d'entre eux.

**Question 5 :** Construire un intervalle de confiance à 95% pour la vraie proportion des casques qui montreraient des dommages à ce test.

Pour répondre à cette question nous avons utilisé la formule 4, pour le code python on a utilisé la fonction *intervalle\_confiance*(*n*, *x*, *ic*) réalisé à la question 3.

```

1 n_casques = 50
2 n_dommages = 18
3 borne_inf_dommages_95, borne_sup_dommages_95 = intervalle_confiance(n_casques,
4                               n_dommages, 95)
5 print("Casques endommagés")

```

```
6 print("\t --> Borne inférieur à 95%:", borne_inf_dommages_95)
7 print("\t --> Borne supérieur à 95%:", borne_sup_dommages_95, end="\n\n")
```

Listing 9 – Code Python question 5

```
1 Casques endommagés
2 --> Borne inférieur à 95%: 0.227
3 --> Borne supérieur à 95%: 0.493
```

Listing 10 – Résultat du code

**Question 6 :** En utilisant un estimateur ponctuel de  $p$  à partir de 50 casques, combien de casques doivent être testés pour avoir une erreur inférieure à 0,02 pour l'IC à 95% de la proportion  $p$  ?

Pour répondre à cette question, nous avons utilisé la formule 5 (cf. Question 4).

```
1 ic = 95
2 alpha = 1 - ic / 100
3 p = 18/50
4 err = .02
5 z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
6 n = (z/err)**2 * p*(1-p)
7 n = np.ceil(n)
8
9 print("Nombre de casque à tester:", n, end='\n\n')
```

Listing 11 – Code Python question 6

```
1 Nombre de casque à tester: 2213.0
```

Listing 12 – Résultat du code

**Question 7 :** Quelle doit être la taille de l'échantillon pour obtenir une erreur inférieure à 0,02 pour l'IC à 95% de la proportion  $p$ . Vous considérerez que vous ne connaissez ni la valeur de  $p$  ni celle de  $\hat{p}$  ? Vous déterminerez la valeur  $p$  pour que la fonction  $n = f(p)$  soit maximale.

En théorie le cas le plus défavorable a lieu quand  $p = 0.5$  (soit  $n$  maximale), mais pour répondre à la question et vérifier cette théorie nous avons décidé de créer un vecteur  $p$ , contenant des valeurs de 0 à 1 par pas de 0,1. En calculant  $n$  avec les différentes valeurs de  $p$ , nous pouvons obtenir la valeur max de  $n$  avec la fonction *max*.

Pour illustrer et vérifier la valeur obtenue, on a tracé la fonction  $n = f(p)$  et la valeur max obtenue, on trouve bien une valeur de  $p = 0.5$  (cf. Figure 1).

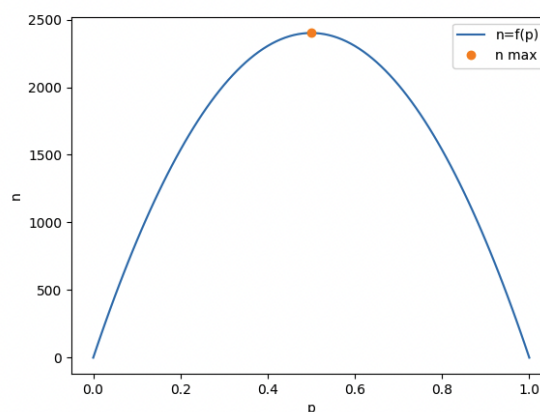


FIGURE 1 – Courbe de  $n = f(p)$

Nous avons utilisé la formule 5 (cf. Question 4).

```
1 ic = 95
2 alpha = 1 - ic / 100
3 err = .02
4 z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
```

```

5 p = np.linspace(0, 1, 101)
6 n = (z/err)**2 * p*(1-p)
7 index_n_max = np.where(n == max(n))
8 n = np.ceil(float(n[index_n_max]))
9
10 print("Question 7:\n", "Nombre de casque à tester:", n, end='\n\n')

```

Listing 13 – Code Python question 7

```

1 Nombre de casque à tester: 2401.0

```

Listing 14 – Résultat du code

## 4 Code complet

```

1 from scipy import stats
2 import numpy as np
3
4 # question 1
5 x = np.array([0.82, 0.87, 0.77, 0.96, 0.75, 0.83, 0.87, 0.81])
6 xn = np.mean(x)
7 ecart_type_x = np.std(x, ddof=1)
8 n = len(x)
9 ic = 95
10 alpha = 1 - ic / 100
11 t = stats.t.ppf(1 - (alpha / 2), n - 1)
12
13 I = xn - (t * ecart_type_x) / np.sqrt(n)
14 u = xn + (t * ecart_type_x) / np.sqrt(n)
15
16 I = round(I, 3)
17 u = round(u, 3)
18
19 print("Question 1:\n", "Borne inférieure:", I, "\n", "Borne supérieure:", u, end="\n\n")
20
21 # question 2
22 x = np.array([0.84, 0.87, 0.89, 0.73, 0.84, 0.81, 0.88, 0.85, 0.89, 0.79, 0.79, 0.90,
23              0.59, 0.75, 0.67, 0.76, 0.86, 0.88, 0.70, 0.75, 0.81, 0.77, 0.83, 0.84,
24              0.71, 0.78, 0.59, 0.91, 0.74, 0.68, 0.77, 0.66, 0.80, 0.74, 1.02, 0.91,
25              0.55, 0.84, 0.66, 0.77])
26 p = np.mean(x)
27 ecart_type_x = np.std(x, ddof=1)
28 n = len(x)
29 ic = 95
30 alpha = 1 - ic / 100
31 z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
32
33 borne_inf = p - z * (ecart_type_x/np.sqrt(n))
34 borne_sup = p + z * (ecart_type_x/np.sqrt(n))
35
36 borne_inf = round(borne_inf, 3)
37 borne_sup = round(borne_sup, 3)
38
39 print("Question 2:\n", "Borne inférieure:", borne_inf, "\n", "Borne supérieure:",
40       borne_sup, end="\n\n")
41
42 # question 3
43 n = 1000
44 n_dupond = 500
45 n_durand = 250
46 n_duroc = 50
47 def intervalle_confiance(n, x, ic):

```

```

48     p = x/n
49     alpha = 1 - ic / 100
50     z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
51     borne_inf = p - z * (np.sqrt(p*(1-p)/n))
52     borne_sup = p + z * (np.sqrt(p*(1-p)/n))
53
54     return round(borne_inf, 3), round(borne_sup, 3)
55
56 borne_inf_dupond_95, borne_sup_dupond_95 = intervalle_confiance(n, n_dupond, 95)
57 borne_inf_durand_95, borne_sup_durand_95 = intervalle_confiance(n, n_durand, 95)
58 borne_inf_duroc_95, borne_sup_duroc_95 = intervalle_confiance(n, n_duroc, 95)
59 borne_inf_dupond_99, borne_sup_dupond_99 = intervalle_confiance(n, n_dupond, 99)
60 borne_inf_durand_99, borne_sup_durand_99 = intervalle_confiance(n, n_durand, 99)
61 borne_inf_duroc_99, borne_sup_duroc_99 = intervalle_confiance(n, n_duroc, 99)
62
63 print("Question 3:")
64 print(" Dupond")
65 print("\t --> Borne inférieur à 95%:", borne_inf_dupond_95)
66 print("\t --> Borne supérieur à 95%:", borne_sup_dupond_95)
67 print("\t\t\t -----")
68 print("\t --> Borne inférieur à 99%:", borne_inf_dupond_99)
69 print("\t --> Borne supérieur à 99%:", borne_sup_dupond_99)
70 print(" Durand")
71 print("\t --> Borne inférieur à 95%:", borne_inf_durand_95)
72 print("\t --> Borne supérieur à 95%:", borne_sup_durand_95)
73 print("\t\t\t -----")
74 print("\t --> Borne inférieur à 99%:", borne_inf_durand_99)
75 print("\t --> Borne supérieur à 99%:", borne_sup_durand_99)
76 print(" Duroc")
77 print("\t --> Borne inférieur à 95%:", borne_inf_duroc_95)
78 print("\t --> Borne supérieur à 95%:", borne_sup_duroc_95)
79 print("\t\t\t -----")
80 print("\t --> Borne inférieur à 99%:", borne_inf_duroc_99)
81 print("\t --> Borne supérieur à 99%:", borne_sup_duroc_99, end="\n\n")
82
83 # question 4
84 ic = 95
85 alpha = 1 - ic / 100
86 p = 17/100
87 err = 1/100
88 z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
89 n = (z/err)**2 * p*(1-p)
90 n = np.ceil(n)
91
92 print("Question 4:\n", "Nombre de personne à interroger:", n, end='\n\n')
93
94 # question 5
95 n_casques = 50
96 n_dommages = 18
97 borne_inf_dommages_95, borne_sup_dommages_95 = intervalle_confiance(n_casques,
98     n_dommages, 95)
99
100 print("Question 5:")
101 print(" Casques endommagés")
102 print("\t --> Borne inférieur à 95%:", borne_inf_dommages_95)
103 print("\t --> Borne supérieur à 95%:", borne_sup_dommages_95, end="\n\n")
104
105 # question 6
106 ic = 95
107 alpha = 1 - ic / 100
108 p = 18/50
109 err = .02
110 z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
111 n = (z/err)**2 * p*(1-p)
112 n = np.ceil(n)

```

```
113 print("Question 6:\n", "Nombre de casque à tester:", n, end='\n\n')
114
115 # question 7
116 ic = 95
117 alpha = 1 - ic / 100
118 err = .02
119 z = stats.norm.ppf(1 - (alpha / 2), loc=0, scale=1)
120 p = np.linspace(0, 1, 101)
121 n = (z/err)**2 * p*(1-p)
122 index_n_max = np.where(n == max(n))
123 n = np.ceil(float(n[index_n_max]))
124
125 print("Question 7:\n", "Nombre de casque à tester:", n, end='\n\n')
```

Listing 15 – Code Python complet TP3