

MAIS OÙ SONT LES NEIGES D'ANTAN ?

OU . . . LES STATISTIQUES

DANS LES NOUVEAUX PROGRAMMES DE COLLÈGE

Jean Claude GIRARD

IREM de LYON

Introduction

L'objet de cet atelier est d'illustrer, sur un exemple, la méthode statistique (ses principes, ses difficultés) et de montrer ce qu'il est possible de faire au collège sur ce sujet. Je commencerai donc par une réflexion sur les nouveaux programmes de collège (en particulier de troisième).

Le projet initial donnait pour objectif "la comparaison de deux séries statistiques concernant un même caractère", non seulement à partir du calcul de paramètres de tendance centrale (moyenne, médiane) ou de position (quartile, décile, quantile) mais aussi de dispersion (initiation à l'écart type). Il mettait l'accent sur la perte d'information liée à l'utilisation de ces paramètres en lieu et place de la série complète et privilégiait une attitude critique vis à vis de ces résumés dans des situations permettant de leur donner du sens et d'en montrer les limites. Chacun pouvait se réjouir de voir l'enseignement des statistiques au collège trouver une cohérence qui pouvait laisser espérer qu'il lui serait enfin donné tout l'intérêt et le temps qu'il mérite.

Malheureusement, la concertation auquel le GTD a soumis le projet a conduit à "un rejet massif" de ce qui a été perçu comme "une inflation concernant les paramètres de dispersion des séries statistiques" (P. Attali).

Alors, que reste-t-il de tout cela dans le programme définitif? Si on le lit un peu vite (et je crains que ce soit le cas le plus fréquent, pour de nombreuses raisons) pas grand-chose! Exit les quantiles, déciles, quartiles et écart type. Plus grave, il n'est plus fait référence à des activités donnant du sens aux calculs statistiques. Il reste seulement une référence à l'éducation du citoyen et la recommandation d'habituer les élèves à avoir une attitude "de lecteur responsable face aux informations de nature statistique". Autrement dit : danger, statistiques!

On peut toutefois faire une relecture plus intéressante de ce programme même si (encore une fois) je crains que ce ne soit pas celle qui sera faite. Malgré une amputation importante, il subsiste deux paragraphes intitulés "caractéristiques de position" et "approche des caractéristiques de dispersion". La médiane a survécu et on a vu

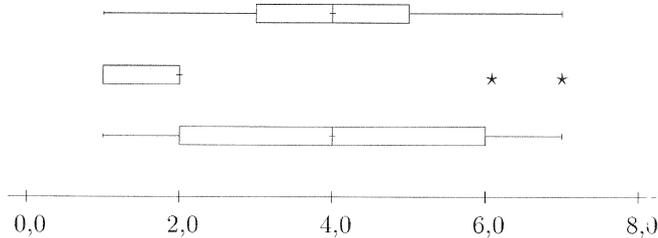
apparaître l'étendue (maximum-minimum) avec la précision "on introduira l'étendue de la série ou de la partie obtenue après élimination de valeurs extrêmes". Cette phrase sibylline conduira de nombreux professeurs à effectuer la soustraction précédente et à conclure que l'idée donnée par cette valeur reflète bien mal la dispersion de l'ensemble des valeurs. Ainsi les trois séries suivantes ont la même étendue :

1	1	*
2	2	**
3	3	***
4	6	*****
5	3	***
6	2	**
7	1	*

1	8	*****
2	8	*****
3	0	
4	0	
5	0	
6	1	*
7	1	*

1	2	**
2	2	**
3	2	**
4	2	**
5	2	**
6	2	**
7	2	**

Par contre, le programme laisse la possibilité d'enlever des valeurs extrêmes (sans dire combien ni pourquoi!). On peut enlever d'autorité 5%, 10% des valeurs à chaque extrémités ou, encore mieux, 25% ce qui laisse la moitié centrale de la série. Le nouveau minimum est alors le premier quartile de la série de départ et le maximum est le troisième quartile. L'écart entre les deux (l'étendue de la nouvelle série) correspond à l'écart inter-quartile de la série de départ ce qui permet de dessiner des boîtes à moustaches¹ pour représenter une série, ou mieux, pour comparer deux séries² et de répondre ainsi à une demande du programme "on pourra ainsi aborder la comparaison de deux séries en calculant quelques caractéristiques de position et de dispersion ou en interprétant des représentations graphiques données". Par exemple, pour les trois séries précédentes :



On peut donc aller plus loin que le calcul d'une moyenne ou d'une médiane ou le tracé d'un histogramme au collège. Encore faut-il que ces calculs ou ces graphiques servent à quelque chose, c'est-à-dire à répondre à une question. C'est ce que je vais essayer de montrer sur un exemple avec l'objectif de développer "votre attitude d'auditeurs responsables vis à vis des statistiques".

1. Pour la construction des graphiques en boîte et un exemple détaillé d'utilisation, voir J.C. Girard, "La médiane, pour quoi faire", *Enseigner la Statistique du CM à la Seconde. Pourquoi? Comment?*, Groupe Probabilités et Statistique, IREM de Lyon, 1998.

2. Le tableau synoptique pour l'ensemble du collège précise en tout et pour tout pour la classe de 3^{ème} ; "Approche de la comparaison de séries statistiques". Ceci semble donc l'objectif de fin de collège.

Mais où sont les neiges d'antan³?

Réflexions générales

Une étude statistique ne fonctionne pas à vide à partir d'une série de nombres sur lesquels on serait amené à faire des calculs ou à faire des graphiques sans savoir pourquoi on les fait ni ce qu'on peut en conclure quand on a terminé.

Une étude statistique est une chaîne qui a pour origine un problème ou une question. Le cœur de l'étude consiste à traiter par des outils mathématiques des mesures obtenues à partir de l'observation d'une variable ou par l'utilisation d'un instrument. Ceci pose le problème des erreurs de mesure ainsi que celui du choix de la variable et de l'instrument pour qu'ils permettent effectivement de répondre à la question posée. Le dernier maillon est constitué par une conclusion finale (peut-être provisoire) ou une décision (qui peut consister à ne rien faire ou à ne pas énoncer d'affirmation péremptoire).

Le problème

La question qui nous servira de point de départ est la suivante : "La planète se réchauffe-t-elle?". La question est d'importance puisque, dans l'affirmative, les simulations ont montré qu'une élévation de la température de 0,5°C aurait des conséquences désastreuses. Il serait bien sûr présomptueux de vouloir trancher alors que les experts sont en désaccord. Notre objectif sera plus modestement d'illustrer une méthode et certains concepts statistiques étudiés au collège.

Choix de la variable

Quelle que soit la variable choisie, elle devra faire l'objet de mesure sur une période assez longue (plusieurs dizaines d'années). Il sera donc nécessaire d'utiliser les résultats d'observations faites par d'autres bien que cela ne soit pas sans danger, car on ne sait pas exactement comment elle ont été obtenues. Il convient donc que ces mesures soient fiables et comparables sur une longue période. C'est difficile mais pas désespéré ! On peut se rappeler, par exemple, que Kepler a découvert les lois qui portent son nom sur le mouvement des planètes⁴ à partir d'une étude statistique des mesures effectuées plusieurs années auparavant par Tycho-Brahé (1546-1601).

La variable qui vient à l'esprit naturellement pour répondre à notre question est la température en un même lieu. Il existe, en effet, des statistiques sur les températures à Paris depuis bien avant la Révolution. Elles permettent, par exemple, de faire apparaître des modifications du climat comme celles qui correspondent au "petit âge glaciaire". Mais les changements intervenus au cours de ce siècle (s'ils s'avèrent se confirmer) ne sont pas de cette amplitude. Une première difficulté serait de déterminer si l'on va privilégier les moyennes mensuelles, les maxima, les

3. F. Villon, *Ballade des dames du temps jadis*.

4. *Astronomia Nova 1609, Harmonices Mundi 1618*.

minima ... Une autre difficulté, bien plus importante, est liée aux instruments de mesure. Les thermomètres du début du siècle n'avaient pas la même précision que les instruments actuels. Ils n'étaient pas construits dans les mêmes matériaux et les méthodes de mesure n'étaient pas standardisées. Une erreur systématique de 1°C (importante) risque alors de passer inaperçue ou au contraire d'être trompeuse.

Pour toutes ces raisons (et aussi parce que les statistiques sont disponibles!) nous avons choisi de travailler sur une autre variable : le nombre de jours de neige annuel⁵. Les observations peuvent être délicates certains jours, la différence entre pluie et neige n'est pas toujours facile, mais ni plus ni moins qu'il y a 100 ans et les mesures ne seront pas entachées d'erreurs systématiques dues aux instruments, l'œil est resté le même depuis le début du siècle.

Le choix de cette variable, ou de toute autre, peut (et doit) être critiqué car de la pertinence (ou non) de cette variable avec la question posée, dépendra la pertinence de la réponse apportée. Aucun calcul mathématique, aussi rigoureux qu'il soit, ne pourra compenser le choix d'une mauvaise variable et d'un mauvais instrument de mesure.

Les données

Pour le tableau de la page suivante, qui donne les nombres de jours de neige à Paris (Parc Montsouris) depuis le début du siècle⁶, l'ordinateur fournit les résultats suivants :

N	MEAN	MEDIAN	TRMEAN	STDEV
98	14.296	13.000	13.955	7.984
MIN	MAX	Q1	Q3	
1.000	36.000	8.000	19.000	

N : Effectif Mean : Moyenne Median : Médiane
 TRMEAN : Moyenne Tronquée à 5% STDEV : Écart Type σ_{n-1}
 MIN : Minimum MAX : Maximum
 Q1 : Premier Quartile Q3 : Troisième Quartile

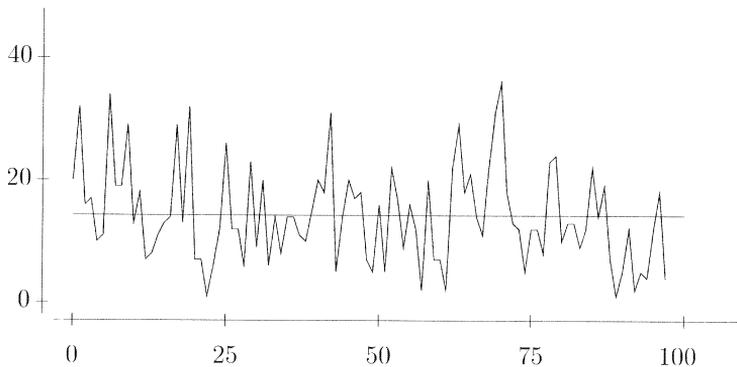
5. En fait la variable qui fait l'objet de l'observation est une variable qualitative à deux modalités : neige ou pas de neige dans la journée, à partir de laquelle on construit la variable quantitative "nombre de jours de neige dans l'année".

6. source QUID 1998, Éditions Robert Laffont

1900	20	1920	7	1940	20	1960	7	1980	10
1901	32	1921	7	1941	18	1961	2	1981	13
1902	16	1922	1	1942	31	1962	22	1982	13
1903	17	1923	6	1943	5	1963	29	1983	9
1904	10	1924	12	1944	14	1964	18	1984	12
1905	11	1925	26	1945	20	1965	21	1985	22
1906	34	1926	12	1946	17	1966	14	1986	14
1907	19	1927	12	1947	18	1967	11	1987	19
1908	19	1928	6	1948	7	1968	22	1988	7
1909	29	1929	23	1949	5	1969	31	1989	1
1910	13	1930	9	1950	16	1970	36	1990	5
1911	18	1931	20	1951	5	1971	18	1991	12
1912	7	1932	6	1952	22	1972	13	1992	2
1913	8	1933	14	1953	17	1973	12	1993	5
1914	11	1934	8	1954	9	1974	5	1994	4
1915	13	1935	14	1955	16	1975	12	1995	12
1916	14	1936	14	1956	12	1976	12	1996	18
1917	29	1937	11	1957	2	1977	8	1997	4
1918	13	1938	10	1958	20	1978	23		
1919	32	1939	15	1959	7	1979	24		

Explication et analyse

Peut-on répondre à notre question à partir de ces données, et si oui, comment? Il est difficile de voir sur le tableau des valeurs de la série complète si le nombre de jours de neige a diminué ou augmenté au cours de ces années, et de combien. Une première idée serait de représenter graphiquement la variable étudiée en fonction de l'année d'observation (on parle dans ce cas de série chronologique).



Ce graphique ne fait pas apparaître d'évidences!

Une autre possibilité est de partager la série initiale de 98 valeurs en deux séries

de 49, puis comparer ces deux nouvelles séries pour voir s'il y a eu des changements perceptibles entre la première moitié et la deuxième moitié de ce siècle.

	N	MEAN	MEDIAN	TRMEAN	STDEV
00-48	49	15.27	14.00	15.02	7.92
49-97	49	13.33	12.00	12.96	8.01
	MIN	MAX	Q1	Q3	
00-48	1.00	34.00	9.50	19.50	
49-97	1.00	36.00	7.00	18.50	

La comparaison des deux séries s'effectue généralement sur des résumés (résultats du calcul de certains paramètres). Le plus classique est la moyenne qui, comme le rappelle le programme officiel, est sensible aux valeurs "aberrantes" ou extrêmes (on verra plus loin s'il y en a et comment on les détermine). De plus, c'est un paramètre de tendance centrale qui ne tient pas compte de la dispersion des valeurs de la série.

Une réponse à la première critique consiste à examiner la médiane⁷ ou la moyenne tronquée. La différence des moyennes est égale à 1,94 est la différence des médianes est à peu près identique puisqu'elle vaut 2 et celle des moyennes tronquées (à 5%) est plus grande puisqu'elle vaut 2,06.

Il semble donc que le nombre de jours de neige dans la première moitié du siècle dépasse de 2 unités celui de la deuxième moitié, en moyenne.

La deuxième critique nous amène à étudier comment les valeurs de la série sont distribuées autour de la moyenne. Une mesure de cette dispersion est l'écart type σ qui est la racine carrée de la variance $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$. Les deux écarts types sont à peu près égaux ici.

La combinaison de l'écart type et de la moyenne permet de donner un meilleur résumé d'une série statistique mais, là encore, l'idée peut être trompeuse. Par exemple, les quatre séries suivantes⁸ ont même moyenne ($m=4$) et même écart type ($\sigma_{n-1}=1$, $\sigma_n=0,961$) :

7. Rappel sur la médiane: la médiane est la valeur qui partage la série en deux séries de même effectif. Autrement dit, si l'on range dans l'ordre croissant les n valeurs de la série, la médiane a pour rang $\frac{n+1}{2}$. Si n est impair, cette valeur est celle du milieu de la série. Si n est pair, on prend la moyenne des deux valeurs du milieu.

8. D'après un article de *Teaching Statistics*.

Histogram of serie 1 N = 13

Midpoint	Count
2	1 *
3	3 ***
4	4 ****
5	5 *****
6	0

Histogram of serie 3 N = 13

Midpoint	Count
2	1 *
3	2 **
4	7 *****
5	2 **
6	1 *

Histogram of serie 2 N = 13

Midpoint	Count
2	0
3	5 *****
4	4 ****
5	3 ***
6	1 *

Histogram of serie 4 N = 13

Midpoint	Count
2	0
3	6 *****
4	1 *
5	6 *****
6	0

Ces séries ne sont manifestement pas du même type. On pourrait les caractériser par la forme de leur distribution (en J, en cloche, en U) ou par d'autres paramètres.

Par exemple, on peut résumer une série par ses trois quartiles⁹. Ils ont l'avantage d'être des paramètres de position et de fournir en même temps une mesure de la dispersion. En effet, l'écart inter-quartile c'est-à-dire la différence entre le premier et le troisième quartile ($Q_3 - Q_1$) est une alternative à l'écart type pour mesurer la dispersion. Moins sensible aux valeurs aberrantes, il est possible, de plus, d'en construire une représentation graphique.

Le calcul des quartiles peut être facilité en représentant la série par un graphique en tiges et feuilles¹⁰, par exemple pour les hauteurs de neige :

Stem-and-leaf of 00-48 N = 49
Leaf Unit = 1.0

1	0	1
12	0	56667777889
(16)	1	0011122233344444
21	1	567788899
12	2	00003
7	2	699
4	3	1224

Stem-and-leaf of 49-97 N = 49
Leaf Unit = 1.0

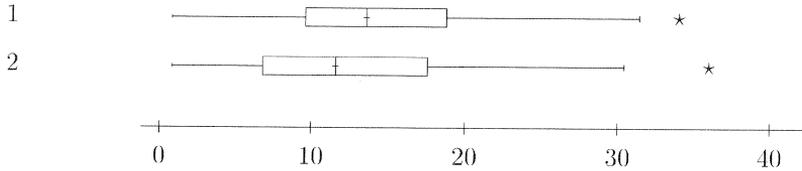
6	0	122244
17	0	55555777899
(14)	1	01222222233344
18	1	6678889
11	2	01222234
3	2	9
2	3	1
1	3	6

On obtient les valeurs des quartiles données par l'ordinateur précédemment.

9. Rappel sur les quartiles : les quartiles sont les valeurs de la série étant classées dans l'ordre croissant de $x_{(1)}$ à $x_{(n)}$, Q_1 représente la valeur de l'observation de rang $\frac{(n+1)}{4}$ et Q_3 l'observation de rang $\frac{3(n+1)}{4}$. Si ces rangs ne sont pas entiers, on effectue une interpolation.

10. Pour une présentation des graphiques en tiges et feuilles, voir par exemple, J.C. Girard. "Des diagrammes à l'histogramme", *Enseigner la Statistique du CM à la seconde. Pourquoi? Comment?*, Groupe des probabilités et Statistique, IREM de Lyon, 1998.

Les écarts inter-quartiles sont alors de 10 et 11,50. On peut les représenter graphiquement sur un graphique en boîte à moustaches (box and whiskers plot).



Deux valeurs sont déclarées aberrantes et notées * parce qu'elles se trouvent à plus de 1,5 fois l'écart inter-quartile de Q_3 .

Les quartiles de la première série dépassent respectivement de 2,5 et 1 ceux de la deuxième série.

Les quartiles forment un résumé d'une distribution. Un résumé perd toujours de l'information mais comparer les deux distributions dans leur entier est difficile. On peut le faire, cependant, en accolant les deux graphiques en tiges et feuilles (surtout si les séries ont le même effectif) :

	1	0	122244
	9887776665	0	55555777899
4444433322211100		1	0122222233344
	9988877765	1	6678889
	30000	2	01222234
	996	2	9
	4221	3	1
		3	6

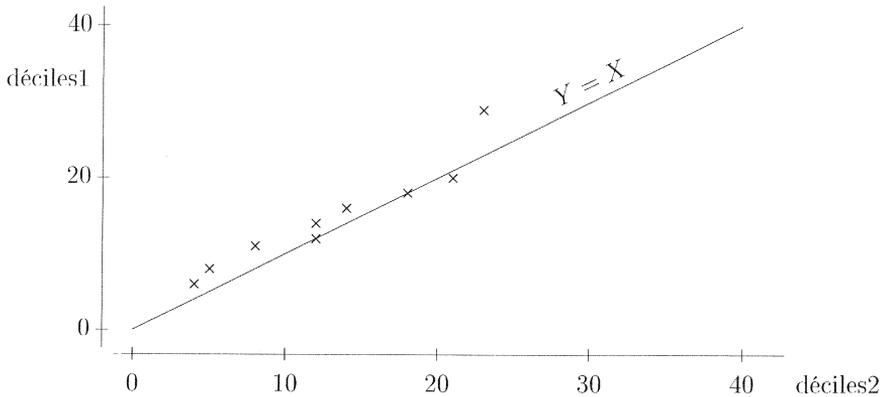
On peut observer que la deuxième série semble être décalée par rapport à la première. Pour des séries d'effectifs différents, on pourrait calculer les déciles¹¹ qui constituent un résumé de 9 nombres de la série.

Les résultats pour les deux séries sont :

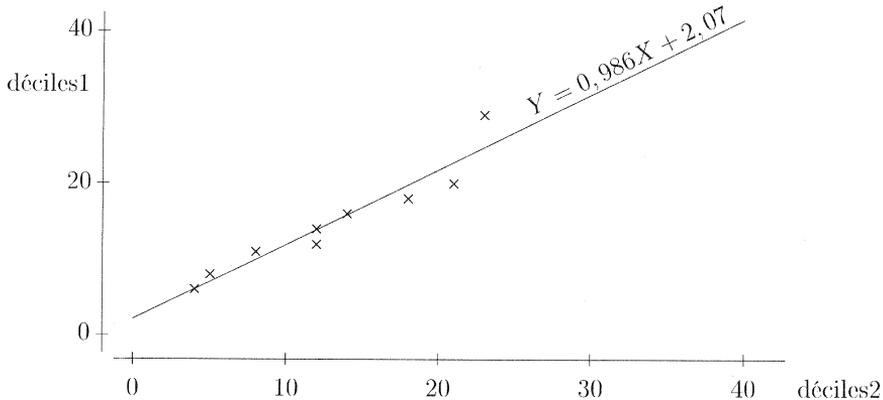
Décile	10%	20 %	30 %	40 %	50 %	60 %	70 %	80 %	90 %
1900	6	8	11	12	14	16	18	20	29
1948									
1949	4	5	8	12	12	14	18	21	23
1997									

Si les deux séries étaient identiques, leurs déciles seraient égaux. Dans un repère, les points dont les coordonnées sont les déciles des deux séries seraient alignés sur la première bissectrice. On constate sur le graphique que ce n'est pas le cas.

11. Rappel du calcul des déciles. Les n valeurs de la série étant classées dans l'ordre croissant de $x_{(1)}$ à $x_{(n)}$, les déciles sont les valeurs des observations dont les rangs sont $\frac{k(n+1)}{10}$, pour $k = 1$ à 9. Si ces rangs ne sont pas entiers, on effectue une interpolation.



On constate plutôt que la droite est décalée d'environ deux unités.



La droite des moindres carrés calculée sur ces 9 points a pour équation $y = 0,986x + 2,07$. Le coefficient de corrélation vaut 0,954.

Régression robuste (et simple)

La régression est bien sûr hors de propos au collège mais on peut (encore?) calculer l'équation d'une droite. Le problème de trouver une droite qui passe le plus près possible de tous les points admet plusieurs réponses suivant la manière de définir ce qu'on entend par "plus près". La droite des moindres carrés est une réponse possible. Une alternative consiste à chercher la droite qui passe par le point médian et qui a pour pente la médiane de toutes les pentes obtenues en prenant les points deux à deux. L'ordinateur calcule ces 36 pentes et leur médiane :

N	N*	MEAN	MEDIAN	TRMEAN	STDEV
35	1	1.095	0.889	0.996	0.732
MIN	MAX	Q1	Q3		
0.250	4.500	0.700	1.211		

La pente est donc 0,889. L'ordonnée à l'origine de cette droite est obtenue par le point dont les coordonnées sont les médianes des deux séries à savoir 14 et 12.

$$14 = 0,889 \times 12 + b$$

$$\text{d'où } b = 3,332$$

La droite a donc pour équation $y = 0,889x + 3,332$.

Si on prend pour pente la moyenne (ou la moyenne tronquée), soit environ 1, on trouve pour équation $Y = X + 2$.

Conclusion de la partie descriptive

Ceci confirme que chaque valeur de rang i de la première série est de 2 à 2,5 unités en dessus de la valeur de même rang de la deuxième série.

Autrement dit on a bien, environ, deux jours de neige de moins par an dans la deuxième moitié du 20^{ème} siècle.

Prolongement inférentiel

La question que l'on se pose maintenant est de savoir si la différence de 2 jours que l'on a observée, en moyenne, entre les deux séries est significative ou si elle pourrait être le fruit du hasard alors que le climat n'a pas changé. En effet, même dans cette dernière hypothèse, on observerait des variations d'une année sur l'autre donc d'un demi-siècle sur l'autre. À partir de quelle différence cette variabilité ne suffit-elle plus à expliquer les différences de moyennes et doit-on remettre en cause l'hypothèse de stabilité du climat? Cette question se traite généralement à un niveau Post-Bac par un test d'égalité des moyennes mais une autre méthode peut être envisagée dès le collège.

Test d'égalité des moyennes

On modélise la situation de la manière suivante : les deux séries sont considérées comme des échantillons de 49 mesures prises sur des années extraites chacune d'un ensemble (une population) d'années. Le test porte sur l'hypothèse (H_0) "les deux échantillons ont été extraits de populations dont les moyennes sont égales" ou "la différence des moyennes est nulle". Dans le cas présent (H_0) signifie que les climats dans les deux moitiés du siècle sont de même température moyenne. Un test classique d'égalité des moyennes donne les résultats suivants :

```
MTB >TwoSample 95.0 '00-48' '49-97';
SUBC> Alternative 0.
TWOSAMPLE T FOR 00-40 VS 49-97
```

	N	MEAN	STDEV	SE MEAN
00-48	49	15.27	7.92	1.1
49-97	49	13.33	8.01	1.1

95 PCT CI FOR MU 00-48 - MU 49-97: (-1.3, 5.1)

TTEST MU 00-48 = MU 49-97 (US NE): T= 1.20 P=0.23 DF= 95

Il donne la probabilité d'observer un tel écart sous (H_0) et l'intervalle de confiance à 95% pour la différence des deux moyennes.

Dans l'hypothèse d'égalité des moyennes, un écart de 1,94 (en plus ou en moins) a une probabilité de 23%, ce qui n'en fait pas un événement extraordinaire. On ne peut donc rejeter cette hypothèse. La différence des moyennes se situe dans l'intervalle [-1,3; 5,1] (au niveau de confiance de 95%). Elle n'est donc pas significative.

La méthode du Bootstrap¹²

Le test précédent requiert des hypothèses qui ne sont pas nécessairement vérifiées ici (indépendance, normalité) et des techniques sophistiquées. Une alternative, envisageable au collège¹³, consiste à répéter plusieurs fois l'opération suivante : construire aléatoirement deux groupes de même effectif dans l'ensemble des 98 valeurs puis calculer la différence de leurs moyennes. Si le résultat dépasse 1,94 (en valeur absolue) avec une fréquence relativement importante, un tel écart ne sera pas révélateur d'une différence entre les deux moitiés du siècle (puisque les groupes ont été faits au hasard, le hasard seul peut l'expliquer). Si, au contraire, cette fréquence est très petite alors le hasard seul ne suffira pas à expliquer un tel écart. Ceci est la base des statistiques inférentielles mais n'est pas simple à comprendre. Il faut d'abord se persuader que les résultats d'un tirage à l'autre ne sont pas les mêmes, certains sont positifs, d'autres sont négatifs, certains sont supérieurs à 1,94 d'autres inférieurs . . . On est donc confronté à une épreuve aléatoire dont on ne connaît pas l'ensemble des possibles et sur laquelle les calculs sont bien délicats¹⁴. Seule l'expérimentation est possible (surtout au collège). On peut, par exemple, écrire chacune des 98 valeurs sur un jeton (ou sur une carte) puis après avoir bien mélangé, constituer au hasard deux groupes de même effectif et enfin calculer la différence de leurs moyennes. Voici les résultats obtenus sur 100 tirages (simulés sur un ordinateur) :

12. Voir, par exemple, Arthur Engel, *Les certitudes du hasard*, Aléas Éditeur, Lyon, 1990.

13. Sans vouloir minimiser les difficultés de compréhension mais avec, au contraire, l'objectif de sensibiliser les élèves à différents concepts qu'ils rencontreront plus tard comme le hasard, la variabilité, les tests d'hypothèse, etc.

14. On pourrait calculer la différence des moyennes dans tous les cas (on dit qu'on fait un test des permutations) il y en a C_{98}^{49} ce qui représente un nombre supérieur à 10^{28}

Stem-and-leaf of C37 N = 100

Leaf Unit = 0.10

1	-4	3
6	-3	98730
15	-2	774432000
33	-1	986666554422211100
(18)	-0	8877776665554333100
49	0	00111111222344444568
29	1	00111112333444668
12	2	002333799
3	3	14
1	4	7

Le nombre de valeurs supérieures à 1,94 (en valeur absolue) est 28. Une simulation sur 1 000 tirages a donné un pourcentage de 23,8. Une telle différence n'est donc pas un événement rare sous (H_0). Par conséquent, on est amené à refuser l'hypothèse qu'il y a une réelle différence entre les deux moitiés du siècle.

Conclusion

Les différences apparues dans l'étude descriptive ne sont pas confirmées par l'étude inférentielle. Il convient donc de se méfier de ce que l'on voit sur les graphiques ou des conclusions tirées hâtivement à partir de quelques calculs trop élémentaires.

D'une façon générale, une conclusion (provisoire, la plupart du temps) ne saurait être apportée sans que des résultats convergents soient obtenus par différentes approches.