

LA PLACE DE LA STATISTIQUE  
DANS LES PROGRAMMES DU SECONDAIRE A LA RENTREE 2012 :  
MISE EN ŒUVRE PEDAGOGIQUE ET FORMATION DES ENSEIGNANTS

**Philippe DUTARTE**

**Résumé – La rénovation des programmes de mathématiques du secondaire aboutit, à la rentrée 2012, avec la mise en place des nouveaux programmes de terminale offrant une place plus importante à la statistique. La mise en œuvre de ces programmes suppose trois points d'appui pédagogiques : une introduction des notions à partir de situations « concrètes » avec une problématique (statistique) identifiée ; une valorisation de la démarche d'investigation, prenant notamment appui sur l'exploitation de logiciels ; une évaluation des compétences acquises dans ce domaine ne se limitant pas à l'aspect classique. Le succès de cette mise en œuvre des programmes dans les classes suppose un effort certain de formation des enseignants.**

### **Introduction**

À la rentrée 2012, se mettent en place les nouveaux programmes de mathématiques de terminales S, ES et STI2D-STL, notamment, achevant ainsi une rénovation des programmes du secondaire (le programme de terminale STMG entre, quant à lui, en application à la rentrée 2013). Ces programmes induisent des changements dans l'enseignement de la statistique, tant quantitatifs que qualitatifs. Ils supposent une mise en œuvre pédagogique spécifique nécessitant un accompagnement des enseignants par une formation adaptée.

### **1. La place de la statistique dans les programmes de mathématiques du secondaire à la rentrée 2012**

Le premier changement visible est d'ordre quantitatif : la statistique et les probabilités occupent une place plus importante que par le passé, avec des contenus nouveaux ou introduits plus précocement. L'enseignement des probabilités est introduit dès la classe de troisième (2008), la loi binomiale apparaît en première (2011) et la loi normale en terminale (2012). Un enseignement de statistique inférentielle est mis en place dès la classe de seconde (2009) avec les notions d'intervalle de fluctuation, d'estimation d'une proportion inconnue ou de prise de décision à partir d'un échantillon. Cet enseignement, reposant d'abord essentiellement sur la simulation, est progressivement formalisé en première et terminale à l'aide des outils probabilistes dont on dispose. Le temps à consacrer à la partie probabilités et statistique est quantifié dans les programmes de terminale. Ainsi, en terminale S, « à titre indicatif, on pourrait consacrer la moitié du temps à l'analyse, l'autre moitié se répartissant équitablement entre géométrie et probabilités-statistique », ou, en terminale ES, « à titre indicatif, on pourrait consacrer environ deux tiers du temps à l'analyse et le reste aux probabilités et à la statistique ». Pour mémoire, le précédent programme de terminale scientifique (2002) donnait, « à titre indicatif », la répartition horaire suivante : « analyse 45%

(environ 14 semaines), géométrie 35% (environ 11 semaines), probabilité et statistique 20% (environ 6 semaines) ».

Le second changement, plus subtil à mettre en œuvre, est d'ordre qualitatif. On passe d'un enseignement, avant les années 2000, de techniques de traitement de données statistiques réalisées par les élèves « à la demande » (calculer une moyenne, représenter un histogramme...), à un enseignement de la statistique au sens de science d'investigation des données, d'aide à la prise de décision (avec estimation du risque) ou à l'estimation. Voici quelques marqueurs de cet « esprit de la statistique » dans les programmes actuels.

### *Classe de troisième (2008)*

L'éducation mathématique rejoint ici [étude des séries statistiques] l'**éducation du citoyen** : prendre l'habitude de s'interroger sur la signification des nombres utilisés, sur l'information apportée par un résumé statistique. De même, c'est pour permettre au citoyen d'**aborder l'incertitude et le hasard dans une perspective rationnelle** que sont introduits les premiers éléments relatifs à la notion de probabilité.

Le travail [en statistique] est conduit aussi souvent que possible **en liaison avec les autres disciplines** dans des **situations** où les données sont exploitables par les élèves. L'utilisation d'un **tableur** permet d'avoir accès à des situations plus riches que celles qui peuvent être traitées « à la main ».

La notion de probabilité est abordée à partir d'**expérimentations** qui permettent d'observer les **fréquences** des issues dans des situations familières (pièces de monnaie, dés, roues de loteries, urnes, etc.). La notion de probabilité est utilisée pour **modéliser** des situations simples de la vie courante.

### *Classe de seconde (2009)*

L'objectif est de faire **réfléchir** les élèves sur des **données réelles**, riches et variées (issues, par exemple, d'un fichier mis à disposition par l'INSEE), synthétiser l'information et proposer des représentations **pertinentes**.

L'objectif est d'amener les élèves à un **questionnement** lors des activités suivantes : l'estimation d'une proportion inconnue à partir d'un échantillon ; la prise de décision à partir d'un échantillon.

### *Classes de premières S et ES (2011)*

L'objectif indiqué dans le programme de seconde est repris :

faire réfléchir les élèves sur des données réelles, riches et variées (issues, par exemple, d'un fichier mis à disposition par l'INSEE).

En statistique inférentielle, on s'appuie sur la loi binomiale pour tester une proportion. L'élève doit être capable

d'**exploiter** l'intervalle de fluctuation à un seuil donné, déterminé à l'aide de la loi binomiale, **pour rejeter ou non une hypothèse** sur une proportion.

L'objectif est

d'amener les élèves à **expérimenter** la notion de « différence significative » par rapport à une valeur attendue et à remarquer que, pour une taille de l'échantillon importante, on conforte les résultats vus en classe de seconde.

### *Classes de terminales S et ES (2011)*

Afin de **traiter les champs de problèmes** associés aux données continues, on introduit les lois de probabilité à densité.

(lois uniforme et normale en ES et S, auxquelles s'ajoute la loi exponentielle en S).

Cette partie [probabilités et statistique] se prête particulièrement à l'étude de **problèmes issus d'autres disciplines**.

Le recours aux représentations graphiques et aux **simulations** est indispensable.

La **problématique de la prise de décision**, déjà rencontrée, est travaillée à nouveau avec l'intervalle de fluctuation asymptotique<sup>1</sup>.

Concernant la partie « estimation » (par intervalle de confiance),

les attendus de ce paragraphe sont modestes et sont à exploiter en lien avec les autres disciplines.

## **2. Mise en œuvre pédagogique**

Il ne suffit certes pas d'énoncer la « loi » (des programmes) pour la voir aussitôt mise en œuvre dans les classes. Les programmes sont accompagnés de documents « ressources » fournissant notamment des exemples de mise en œuvre. Les manuels scolaires, édités lors de l'application d'un nouveau programme, jouent un rôle essentiel, demeurant la source documentaire prépondérante des professeurs. L'évaluation, en particulier les examens, détermine en partie les contenus enseignés et leur mode d'enseignement (on parle de « pilotage par l'examen »). Enfin, la formation, initiale et continue, des professeurs doit prendre en compte les objectifs spécifiques de l'enseignement de la statistique.

Nous nous penchons dans cette partie sur ce qui nous semble constituer trois points d'appui pédagogiques incontournables.

### **2.1. En situation « concrète »**

Est-ce utile de le rappeler, la statistique est une science « appliquée », c'est ce qui fait son charme. Cela ne signifie pas que les mathématiques mises en jeu, notamment des résultats de probabilités ou d'analyse, soient de seconde catégorie mais que la compréhension des concepts et des démarches est liée à des usages extérieurs aux mathématiques. Un enseignement de la statistique de type « académique », partant de définitions posées a priori et énonçant des théorèmes en dehors de tout contexte est donc non seulement inefficace, mais néfaste. À juste titre pourrait-on dire que, dans ce cas, mieux vaut ne pas confier l'enseignement de la statistique à un professeur de mathématiques. A contrario, le professeur de mathématiques est le seul capable de

---

<sup>1</sup> On peut regretter que cette remarque, qui nous paraît essentielle, n'apparaisse qu'à la fin des « commentaires » concernant la partie « intervalle de fluctuation ».

détacher les concepts mathématiques mis en œuvre des différents contextes spécifiques, pour mettre en évidence leur caractère universel permettant d’aller au-delà d’un agrégat de recettes ad hoc et donc de s’adapter à des situations imprévues. Il est ainsi essentiel que l’enseignement de la statistique soit pris en charge par un professeur de mathématiques.

Certains types d’exercices nous semblent ne pas avoir leur place dans un apprentissage de la statistique. Il s’agit d’exercices stériles, vides de sens et ne répondant à aucune problématique.

Voici deux exemples (ils semblent se raréfier dans les manuels).

**27** \* On donne la série statistique suivante :

<b>Valeurs</b>	1	a	5	2a	9
<b>Effectifs</b>	28	22	14	20	16

a) Déterminez a sachant que la moyenne est égale à 4,59.  
 b) Calculez la médiane de cette série.

Figure 1 – Manuel de seconde (2009)

**54** [6 points]

Soit  $X$  une variable aléatoire suivant une loi  $\mathcal{N}(0 ; 1)$ .

- Vérifier que les intervalles  $]-\infty ; 1,65[$  et  $[-2,57 ; 1,69]$  peuvent être considérés comme des intervalles de fluctuation de  $X$  au seuil de 95 % (c’est-à-dire des intervalles  $I$  tels que  $P(X \in I) \geq 0,95$ ).
- Montrer qu’il existe une valeur  $a$  minimale telle que l’intervalle  $[-a ; a]$  soit un intervalle de fluctuation de  $X$  au seuil de 95 %. En donner une valeur approchée à  $10^{-2}$  près.
- a. Montrer qu’il existe un unique réel  $b$  tel que  $(-2 \leq X \leq -2 + b) = 0,95$ .  
 b. Prouver que  $b < a + 2$  où  $a$  est la valeur de la question 2.  
 c. Déterminer une valeur approchée de  $b$  à  $10^{-2}$  près.
- Montrer qu’il n’existe aucun réel  $c$  tel que  $P(-1 \leq X \leq -1 + c) = 0,95$ .

Figure 2 – Manuel de terminale S (2012)

Gourmands en temps (il faut faire plusieurs exercices du même type, pour parvenir, si l’on est résistant, à acquérir une compétence inutile, sauf le jour de « l’interro »), ils éloignent de l’apprentissage des capacités attendues, rendant plus difficile la compréhension du sens de notions liées à la modélisation : une moyenne, une médiane, un intervalle de fluctuation correspondent à des calculs plus ou moins arbitraires, selon un degré d’approximation qui l’est tout autant<sup>2</sup>.

<sup>2</sup> On peut parler à ce propos de « réduction arithmétique » ou de « réduction mathématique » de la statistique (voir [3]).

La contextualisation est particulièrement nécessaire en statistique inférentielle. L'objectif est de construire le sens critique, par le choix des méthodes, des paramètres, l'évaluation des risques, les limites du modèle, la signification de la « confiance ». Citons Norbert Meusnier<sup>3</sup> :

L'éducation à l'aléatoire [...] devrait avoir pour but fondamental la prise de conscience que toute décision s'accompagne d'un risque, mais que ce risque peut être évalué.

La contextualisation d'une méthode statistique n'est pas toujours simple à réaliser (et les professeurs doivent être accompagnés) mais l'enjeu est d'importance. Voyons un exemple d'exercice qui, sans être mathématiquement faux, est, du point de vue de la modélisation, maladroitement posé.

**40** Dans tout cet exercice, la production est supposée suffisamment importante pour que l'on assimile le choix d'un échantillon à un tirage avec remise.

Un sous-traitant est chargé de concevoir des pièces pour un constructeur automobile.

**1.** Un sondage est réalisé pour tester la qualité de la production du sous-traitant. On suppose que la proportion de pièces non conformes dans la production est comprise entre 0,02 et 0,98 ; sur un échantillon de 250 pièces, 12 ne répondent pas au cahier des charges.

**a)** Déterminer l'intervalle de confiance au niveau de confiance 95 % obtenu à partir de cet échantillon. Donner les bornes de l'intervalle arrondies à  $10^{-3}$  près.

**b)** Le client veut un taux de pièces non conformes inférieur ou égal à 6 %. Compte tenu du sondage effectué, peut-on affirmer avec un niveau de confiance de 95 % que cette condition est respectée ?

**2.** Dans cette question, toute trace de recherche ou d'initiative sera prise en compte dans la notation.

Le sous-traitant procède à de nouveaux réglages et fait un nouveau sondage sur 250 pièces. Déterminer, à l'aide de la calculatrice, le nombre maximum de pièces non conformes dans cet échantillon pour pouvoir affirmer que, au niveau de confiance de 95 %, le taux de pièces non conformes est inférieur ou égal à 6 % dans la production. Expliquer votre démarche.

Figure 3 - Manuel de terminale STI2D (2012)

Le contexte est celui d'une prise d'échantillon pour un contrôle de qualité, mais la problématique n'apparaît qu'à la question 1.b) : « le client veut un taux de pièces non conformes inférieur ou égal à 6% ». Passons outre le fait qu'un taux de 6% est sans doute peu réaliste pour la survie économique de l'entreprise (on ne prétend pas, en terminale, envisager une situation dans toute la complexité des applications réelles). Ce

---

<sup>3</sup> Voir [4].

qui importe, c'est de comprendre que, dans ce type de contexte, on a sans doute intérêt à effectuer un test (raisonner en termes d'intervalle de fluctuation et de prise de décision, pour reprendre la terminologie du secondaire) puisque l'on a une idée de ce que devrait être la proportion  $p$  de pièces défectueuses dans la production. Notons également que la situation est ici « unilatérale » : on ne voudrait pas que  $p$  dépasse 6% (supposer que la proportion de pièces non conformes dans la production est d'au moins 2% est assez étrange). Ce qui rend un intervalle de confiance à 95% (bilatéral) a priori inadapté.

Bien entendu, un « bon » élève (la question 2. n'est pas facile) peut « faire » cet exercice de mathématiques. On attend à la question 2. la recherche du plus petit entier  $k$  tel que :

$$\frac{k}{250} + 1,96 \frac{\sqrt{\frac{k}{250} \left(1 - \frac{k}{250}\right)}}{\sqrt{250}} < 0,06$$

(borne supérieure de l'intervalle de confiance « standard » à 95% inférieure à 0,06). La calculatrice ou le tableur fournit  $k = 9$ . Mais la démarche est assez compliquée et la confiance de 95% est difficile à interpréter.

En posant la problématique dès le début de l'exercice, un élève de terminale (et même de première) peut procéder autrement (et sans doute mieux). Supposons que, situation à la limite de l'acceptable, la proportion de pièces défectueuses dans la production est  $p = 0,06$ . Sous cette hypothèse, la variable aléatoire  $X$  correspondant au nombre de pièces défectueuses observées sur un échantillon aléatoire de taille 250 suit la loi binomiale de paramètres 250 et 0,06. Les outils de calcul (calculatrice, tableur, GeoGebra4...) fournissent :  $P(X \leq 8) \approx 0,033$  et  $P(X \leq 9) \approx 0,064$  (on a là en quelque sorte des intervalles de fluctuation unilatéraux pour la variable aléatoire  $X$ ). On peut être conduit à la règle de décision suivante : on accepte dans l'échantillon prélevé un maximum de 8 pièces défectueuses. La probabilité d'accepter à tort la production est de 3,3 % (risque de rejet à tort de l'hypothèse  $p = 0,06$  à gauche, c'est-à-dire que dans 3,3% des cas, on accepte une production pour laquelle  $p = 0,06$  ou davantage). On remarque qu'avec la règle fournie par l'intervalle de confiance le risque d'accepter à tort la production est de 6,4%.

Cette méthode (utilisation directe de la loi binomiale) apparaît plus simple et plus précise que celle de l'intervalle de confiance, en un mot mieux adaptée.

#### *a) Prise de décision à l'aide d'un intervalle de fluctuation*

La mise en œuvre d'un intervalle de fluctuation, pour la prise de décision, devrait répondre, au lycée, à un certain nombre de critères, garantissant le respect de l'esprit de la démarche statistique, susceptible, selon les orientations choisies, d'être poursuivie dans l'enseignement supérieur<sup>4</sup>.

- Le premier critère nous semble devoir être de poser un problème, une situation, ayant quelque écho avec le monde réel. Cette situation doit être telle que la proportion  $p$  dans la population est supposée connue : c'est l'hypothèse de départ, avec laquelle on construit la règle de décision. Cela exclut les situations où l'on

---

<sup>4</sup> Satisfaire tous ces critères dans un exercice scolaire de lycée n'est pas évident. Le premier (une situation-problème) est cependant incontournable : le choix d'une méthode statistique dépend de la situation.

n'a aucun a priori sur la valeur de  $p$  (type « sondage »), pour lesquelles on aura recourt à un intervalle de confiance<sup>5</sup>.

- Il faut ensuite que cette situation soit de type « bilatérale » : par rapport à la valeur de  $p$  « visée » (celle de l'hypothèse) on s'intéresse à un écart trop important à gauche et à droite.

- La règle de décision doit être élaborée, autant que possible<sup>6</sup>, avant la prise d'échantillon. Il est plus honnête de décider d'une règle avant de jouer, qu'après la partie<sup>7</sup>. Par ailleurs, l'approche de Neyman-Pearson est fréquentiste. Il ne s'agit pas d'élaborer une règle de décision à partir d'un échantillon, mais d'élaborer une règle de décision en amont dont on sait évaluer les risques sur un grand nombre d'échantillons.

- L'échantillon doit être prélevé par tirage au hasard avec remise (équiprobabilité garantie par randomisation).

- En cas de rejet de l'hypothèse, il faut savoir que l'on peut interpréter le 5% (ou le 1% en terminale S dans le cas d'un intervalle de fluctuation asymptotique à 99%<sup>8</sup>) comme la probabilité de commettre une erreur de décision (probabilité de rejet de l'hypothèse sachant qu'elle est vraie). En cas d'acceptation de l'hypothèse, il faut savoir que l'estimation de la probabilité d'erreur de décision est plus compliquée, puisqu'elle dépend de la valeur alternative de  $p$ <sup>9</sup>. Il faudrait par exemple envisager une autre valeur possible de  $p$  ou balayer différentes valeurs possibles, avec un tableur par exemple.

### Un exemple dans une situation bilatérale

L'exercice suivant, où il est question de variable aléatoire, peut être proposé au niveau première ou terminale.

On considère la couleur des yeux des mouches drosophiles. Par croisement de mouches homozygotes « yeux rouges » de gènes AA et de mouches homozygotes « yeux bruns » aa, on obtient des mouches hétérozygotes Aa. Si l'on croise ces hétérozygotes entre eux,

---

5 L'estimation par intervalle de confiance est une démarche inductive à l'état d'esprit très différent de celui d'un test. On n'y retrouve pas la démarche déductive de type « raisonnement par l'absurde » à l'œuvre dans la prise de décision correspondant aux tests, en cas de rejet de l'hypothèse. Il faut absolument éviter de mener les deux démarches pour une même situation.

6 Dans certains cas, voir par exemple les cas de leucémies à Woburn (document ressources pour le lycée professionnel), on ne peut procéder ainsi. Le raisonnement consiste à supposer que l'échantillon étudié résulte d'un échantillonnage aléatoire sous une certaine hypothèse.

7 Dans nombre d'exercices présents dans les manuels de terminale (2012), on débute par une observation sur un échantillon, puis on envisage une, voire plusieurs règles de décision.

8 Le choix du seuil de décision (1% ou 5%) doit, bien entendu, se faire avant la prise d'échantillon. Il dépend du risque (de première espèce) consenti (si l'on diminue le risque de première espèce, on augmente le risque de seconde espèce, à taille d'échantillon égale). Ce choix dépend des utilisateurs et non du statisticien.

9 Les notions d'erreur de première et de seconde espèces ne sont pas au programme au lycée (aucune autonomie des élèves à ce sujet n'est attendue, ni même la connaissance de ce vocabulaire). Pour autant, il nous semble inconcevable de passer cette problématique (il y a deux types d'erreurs et donc de risques) complètement sous silence, tant elle est essentielle. Une analogie simple suffit à faire comprendre la situation. Une prise de décision est comme un jugement au tribunal. L'hypothèse est que le prévenu est présumé innocent. Il y a deux risques au jugement : celui de condamner un innocent (rejet à tort de l'hypothèse, première espèce), ou d'absoudre un coupable (acceptation à tort de l'hypothèse, seconde espèce).

on doit obtenir, selon les lois de Mendel, dans cette seconde génération, 75% de « yeux rouges » (AA et Aa) et 25% de « yeux bruns » (aa).

On souhaite tester l'hypothèse selon laquelle la proportion de drosophiles « yeux rouges » de seconde génération est  $p = 0,75$  en mettant en place une expérimentation permettant d'observer 300 drosophiles de seconde génération (considérées comme un échantillon aléatoire).

1. Sous l'hypothèse  $p = 0,75$ , déterminer l'intervalle de fluctuation au seuil de 95% de la variable aléatoire correspondant à la fréquence du caractère « yeux rouges » sur un échantillon aléatoire de taille 300.
2. Énoncer la règle de décision permettant de rejeter, ou non, l'hypothèse  $p = 0,75$ , au seuil de 5%, sur un échantillon aléatoire de taille 300.
3. L'expérience permet d'observer 237 « yeux rouges » et 63 « yeux bruns ». Cette répartition est-elle conforme à la loi de Mendel, au seuil de décision de 5% ?

Détaillons le raisonnement que l'on peut suivre.

1. On considère que les 300 drosophiles qui seront observées résultent d'un tirage au hasard et avec remise dans la population des drosophiles de seconde génération. Dans ce cadre, la variable aléatoire  $X$  correspondant au nombre de drosophiles présentant le caractère « yeux rouges » dans un tel échantillon suit, sous l'hypothèse  $p = 0,75$ , la loi binomiale de paramètres  $n = 300$  et  $p = 0,75$ . Il est donc possible, sous l'hypothèse  $p = 0,75$ , de prévoir la « variabilité » de la variable aléatoire  $F = \frac{X}{300}$  correspondant à la fréquence du caractère « yeux rouges » sur un tel échantillon. Compte tenu de la nature du problème (on accepte l'hypothèse  $p = 0,75$  à condition que la fréquence observée ne s'en éloigne pas trop, ni à gauche, ni à droite), on recherche un intervalle de fluctuation « bilatéral » (c'est-à-dire avec une zone de rejet à gauche et une zone de rejet à droite) correspondant à une probabilité d'au moins 95%. Ceci de sorte que le risque<sup>10</sup> de rejeter à tort cette hypothèse soit d'au plus 5% et même, ici, d'au plus 2,5% de chaque côté. On est donc amené à rechercher les entiers  $a$  et  $b$  tels que :  $P(X < a) \leq 0,025$  et  $P(X > b) \leq 0,025$  (et donc  $P(a \leq X \leq b) \geq 0,95$ ).

Le logiciel GeoGebra, dans sa version 4, permet de déterminer aisément (et de manière visuelle) ces valeurs  $a$  et  $b$  pour lesquelles les probabilités « à gauche » et « à droite » dépassent le seuil de 0,025. On trouve  $a = 210$  et  $b = 239$ .

---

<sup>10</sup> L'apport de Pearson et Neyman nous invite à porter le regard sur le risque (ici de première espèce) de sorte que l'intervalle de fluctuation utile à une prise de décision se détermine surtout par considération de son complémentaire. N'importe quel intervalle de fluctuation ne fait pas l'affaire et, d'une certaine façon, il y a ici unicité de la réponse à apporter dans la recherche de l'intervalle de fluctuation adapté au problème.



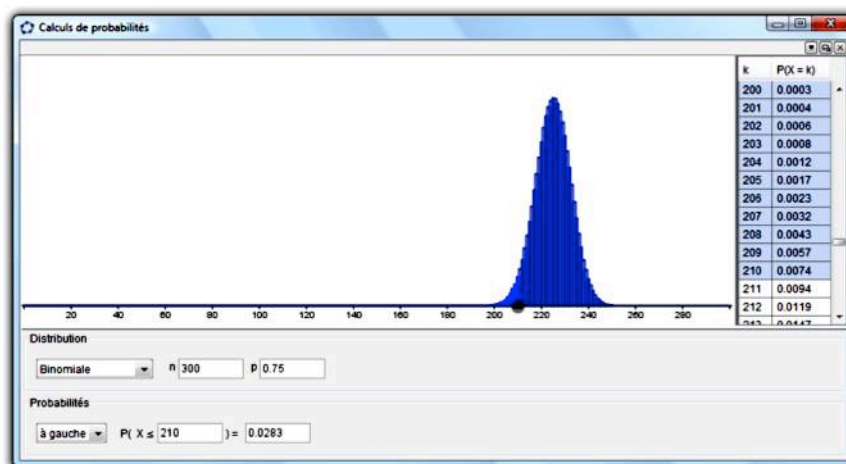


Figure 4 – Avec Geogebra, version 4, on trouve  $a = 210$  et  $b = 239$ .

En divisant par 300, on obtient l'intervalle de fluctuation au seuil de 95% de  $F$  :  $I = [0,7 ; 0,797]$ .

En classe de terminale, on peut préférer utiliser l'intervalle de fluctuation asymptotique, déterminé à partir de la loi normale, et qui fournit un résultat extrêmement proche :  $[0,701 ; 0,799]$ . L'intervalle de fluctuation donné en seconde, majoration du précédent, est  $[0,69 ; 0,81]$ .

2. Après prélèvement d'un échantillon aléatoire de taille 300 et calcul de la fréquence  $f$  du caractère « yeux rouges » sur cet échantillon, la règle de décision, au seuil de 5%, est la suivante : si  $f$  n'appartient pas à l'intervalle  $I$  de fluctuation à 95% précédent, on rejette l'hypothèse  $p = 0,75$ , si  $f$  appartient à  $I$ , on ne rejette pas cette hypothèse.

Il est important de savoir que les 5%<sup>11</sup> correspondent à la probabilité de rejeter à tort l'hypothèse  $p = 0,75$  (c'est comme cela qu'a été construite la procédure de décision).

3. On observe  $f = 0,79$ . Cette valeur appartient à l'intervalle  $I$  de fluctuation à 95%. Cette observation est donc « conforme à » (ou « compatible avec ») l'hypothèse  $p = 0,75$  au seuil de décision de 5%. On s'exprime ainsi car, dans cette situation, il serait un peu étrange d'affirmer qu'on « ne rejette pas la loi de Mendel ». En cas de rejet de l'hypothèse, on serait sans doute amené à étudier les conditions du protocole expérimental.

### Un exemple dans une situation unilatérale

Concernant une proportion, nombre de situations de prise de décision sont, de manière plus « naturelle », unilatérales plutôt que bilatérales. Il nous semble dommage de s'interdire l'accès à ces situations, tout en étant conscient qu'aucune autonomie des élèves n'est attendue à cet égard. Il ne s'agit pas pour autant de taire le problème et de faire, dans une situation annoncée comme unilatérale, « comme si de rien n'était »<sup>12</sup>.

11 Avec la loi binomiale, ce n'est pas exactement 5% (on peut déterminer la probabilité correspondant à la zone de rejet).

12 C'est le cas dans de nombreux exercices présents dans des manuels de terminale (2012) où, dans une situation unilatérale, on semble prendre une décision « au seuil de 95% » alors qu'elle l'est au seuil de 97,5%.

Nous reprenons ici un exemple présenté dans le document ressources pour la classe de terminale<sup>13</sup>, en insérant la question 2., essentielle à la compréhension du « 95% ».

On admet que dans la population d'enfants de 11 à 14 ans d'un département français le pourcentage d'enfants ayant déjà eu une crise d'asthme dans leur vie est de 13%. Un médecin d'une ville de ce département est surpris du nombre important d'enfants le consultant ayant des crises d'asthme et en informe les services sanitaires. Ceux-ci décident d'entreprendre une étude et d'évaluer la proportion d'enfants de 11 à 14 ans ayant déjà eu des crises d'asthme.

Ils sélectionnent de manière aléatoire 100 jeunes de 11 à 14 ans de la ville.

La règle de décision prise est la suivante : si la proportion observée est supérieure à la borne supérieure de l'intervalle de fluctuation asymptotique au seuil de 95% alors une investigation plus complète sera mise en place afin de rechercher les facteurs de risque pouvant expliquer cette proportion élevée.

1. Déterminer l'intervalle de fluctuation asymptotique au seuil de 95% de la proportion de jeunes de 11 à 14 ans ayant eu une crise d'asthme dans un échantillon de taille 100.
2. Indiquer une valeur approchée de la probabilité de mener une investigation supplémentaire à tort (c'est-à-dire alors que la proportion d'enfants de 11 à 14 ans de la ville V ayant déjà eu une crise d'asthme dans leur vie est de 13%) ?
3. L'étude réalisée auprès des 100 personnes a dénombré 19 jeunes ayant déjà eu des crises d'asthme. Que pouvez-vous conclure ?

1. On trouve  $[0,06 ; 0,20]$ . Remarquons qu'il faut lire attentivement l'énoncé pour comprendre qu'il y a deux populations. Celle du département, où  $p = 0,13$  est connu, et celle de la ville V, où l'on fait en quelque sorte, l'hypothèse que  $p = 0,13$ .

2. La problématique de cet exercice est unilatérale, ce qui est clairement affirmé dans la règle de décision. En revanche, l'intervalle de fluctuation asymptotique au seuil de 95% est bilatéral (c'est le seul au programme de terminale). On ne peut donc répondre « 5% » à la question posée. En revanche, un élève de terminale connaît la propriété de symétrie de la courbe de Gauss et sait donc que ces 5% (correspondant aux fréquences observées situées en dehors de l'intervalle de fluctuation lorsque  $p = 0,13$ ) se partagent en 2,5% de chaque côté de l'intervalle de fluctuation. La probabilité demandée est donc d'environ<sup>14</sup> 2,5%.

Il est essentiel, et simple, de comprendre que le risque consenti de mener une investigation supplémentaire inutile (et éventuellement coûteuse) est de 2,5%.

3. La valeur 0,19 est à l'intérieur de l'intervalle de fluctuation asymptotique au seuil de 95%, on en conclut que la règle de décision choisie ne prévoit pas de réaliser une enquête supplémentaire<sup>15</sup>.

---

13 Voir [8] page 22.

14 Le « environ » provenant de l'aspect asymptotique de l'intervalle de fluctuation. Pour un résultat « exact », travailler avec la loi binomiale.

15 On peut faire remarquer que dans cette situation d'acceptation de l'hypothèse, on ne sait pas, en l'état, quantifier le risque. Cela peut-être prétexte à prolongement de l'exercice. Supposons par exemple que la proportion d'enfants de cet âge asthmatiques dans la ville V (connue comme souffrant de pollution de l'air) soit en fait de 25%, quelle est la probabilité que la règle de décision précédente conduite à accepter (à tort) l'hypothèse  $p = 0,13$  ? En terminale, il s'agit d'une question « ouverte » : on peut effectuer des simulations, utiliser la loi binomiale ou calculer  $P(F \leq 0,20)$ , où  $F$  suit la loi normale de moyenne 0,25 et d'écart type .

b) Explorer un grand nombre de données

Le défi de la statistique du XXI<sup>e</sup> siècle est d'exploiter des quantités toujours plus gigantesques de données. Il est essentiel de placer les élèves dans des situations où il s'agit d'explorer un grand nombre de données et en particulier des données « réelles ». C'est l'occasion de mettre les élèves « en situation » d'un usage très actuel de la statistique, de développer des compétences de prise d'initiative, de choix de résumés numériques ou graphiques pertinents, d'outils de modélisation plus ou moins adaptés aux situations, et de montrer l'intérêt et la puissance des méthodes mises en œuvre. Nous évoquons ici deux exemples de traitements de données accessibles sur Internet<sup>16</sup>.

Tailles père-fils (données de Karl Pearson)

Le statisticien britannique Karl Pearson (1857-1936), dans le cadre de recherches sur l'hérédité, a établi un fichier de 1 078 couples de mesures de la taille du père et du fils (adulte). Le fichier tableur fournit ces données, exprimées en mètres.

1. Regrouper les tailles des pères en classes d'amplitude 0,01 mètre. Peut-on ajuster à l'histogramme normalisé des fréquences, correspondant à ces classes, une loi à densité figurant au programme de terminale ?
2. Même question pour les tailles des fils.
3. La taille moyenne d'un homme d'âge compris entre 20 et 29 ans est actuellement, en France, de 1,77 m. En utilisant le modèle de la question précédente estimer la probabilité qu'un jeune homme anglais de ces âges ait une taille supérieure ou égale à 1,77 m à la fin du XIX<sup>e</sup> siècle.
4. Représenter les points de coordonnées  $(x, y)$  correspondant aux tailles (père, fils) (« nuage de points ») et figurer la droite  $D$  d'équation  $y = x$ .
  - a. A-t-on autant de points au-dessus et en dessous de la droite  $D$  ? Donner une interprétation.
  - b. Pour  $x \leq 1,58$ , tous les points sont au-dessus de la droite  $D$  et pour  $x \geq 1,88$ , tous les points sont en dessous de  $D$ . Donner une interprétation.

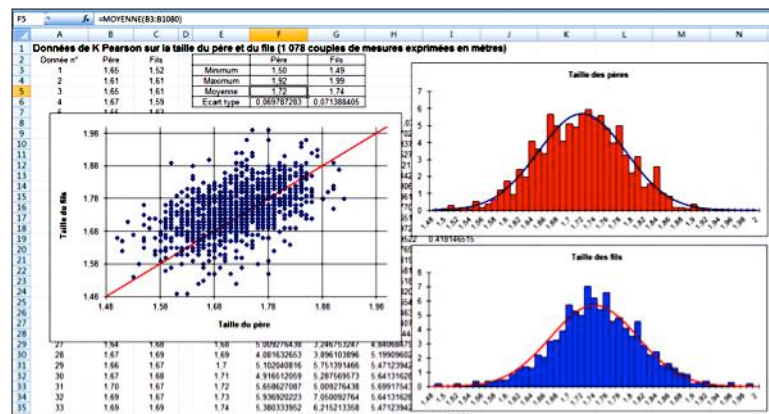


Figure 5 – Tailles père-fils

16 Il suffit, s'agissant ici de données anglo-saxonnes, d'entrer dans un moteur de recherche, *Karl Pearson height data set* ou *Old Faithful data set*.

### Éruptions du geyser Old Faithful

Le geyser Old Faithful est situé dans le parc Yellowstone aux États-Unis. Son nom signifie « vieux fidèle » en raison de la régularité de ses éruptions. Les données statistiques permettent d'étudier cette « fidélité ». Le fichier tableur fournit la durée entre le début de chaque éruption, exprimée en minutes, pour les 5 699 éruptions de l'année 2010, ainsi que la durée moyenne journalière entre éruptions pour chacun des 365 jours de l'année.

1. Regrouper les 5 699 durées entre éruptions en classes d'amplitude 4 minutes. Peut-on ajuster à l'histogramme normalisé des fréquences, correspondant à ces classes, une loi à densité figurant au programme de terminale ?
2. Regrouper les 365 durées moyennes journalières entre éruptions en classes d'amplitude une minute. Peut-on ajuster à l'histogramme normalisé des fréquences, correspondant à ces classes, une loi à densité figurant au programme de terminale ?

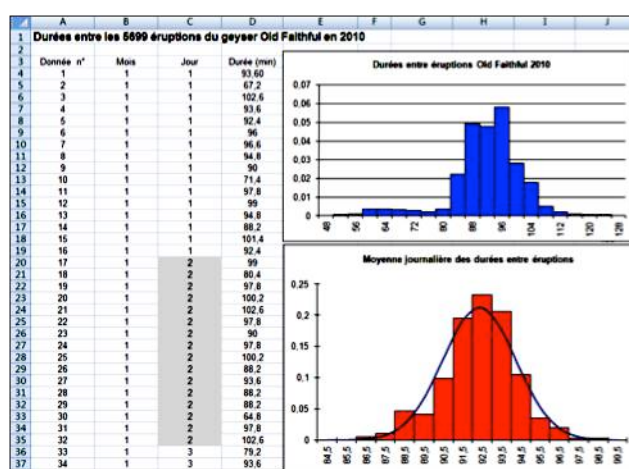


Figure 6 – Eruptions du geyser Old Faithful

### 2.2. En valorisant la démarche d'investigation et l'exploitation de logiciels

Ainsi que le signale l'introduction commune aux disciplines scientifiques des programmes de collège, « dans la continuité de l'école primaire, les programmes du collège privilégient pour les disciplines scientifiques et la technologie une démarche d'investigation », avec des spécificités en mathématiques, en particulier concernant la validation de conjectures à l'aide de démonstrations. Cette approche se poursuit au lycée et les exemples illustratifs seront ici choisis au niveau de la terminale scientifique, où se situe la nouveauté à la rentrée 2012.

Cette démarche « privilégie la construction du savoir par l'élève » : dans le cadre d'une « situation-problème », on réalise des « expériences », pouvant correspondre en mathématiques à une expérimentation à l'aide de logiciels, conduisant notamment à formuler des conjectures ou à mieux comprendre un concept. La simulation joue un rôle particulier dans l'expérimentation, notamment dans l'étude de phénomènes aléatoires.

Nous prenons trois exemples dans le domaine « probabilités-statistique » du nouveau programme de terminale scientifique (2012) montrant le rôle de l'expérimentation à l'aide de logiciels.

### Introduction du théorème de Moivre-Laplace

En terminale S, le théorème de Moivre-Laplace peut être énoncé ainsi (il est admis) :

Soit, pour tout entier  $n$ , une variable aléatoire  $X_n$  qui suit la loi binomiale  $B(n, p)$  et soit  $Z_n = \frac{X_n - np}{\sqrt{np(1-p)}}$ , la variable centrée réduite associée à  $X_n$ .

Alors, pour tous réels  $a$  et  $b$ , tels que  $a < b$  :  $\lim_{n \rightarrow +\infty} P(a \leq Z_n \leq b) = \int_a^b \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx$

Il n'est pas question d'énoncer ce résultat sans expérimentation préalable par les élèves (du moins souhaitons-le) et les commentaires du programme conseillent de s'appuyer sur l'observation graphique.

On peut proposer la démarche suivante (avec ici une présentation relativement « guidée »).

Dans une population, la proportion des personnes possédant le gène A actif est  $p = 0,4$ . On prélève au hasard un échantillon de taille  $n$  dans cette population (celle-ci étant suffisamment grande pour considérer qu'il s'agit de tirages avec remise). On s'intéresse à la répartition de ce gène sur les différents échantillons possibles.

#### A. Histogramme normalisé de la loi binomiale centrée réduite

On désigne par  $X_n$  la variable aléatoire associant à chaque échantillon de taille  $n$  le nombre de personnes de l'échantillon possédant le gène A actif.

1. a. Justifier que la variable aléatoire  $X_n$  suit une loi binomiale dont on donnera les paramètres.

b. Déterminer, en fonction de  $n$ , l'espérance  $m$  et l'écart-type  $\sigma$  de la variable aléatoire  $X_n$ .

c. On considère l'événement  $E_n$  : «  $m - \sigma \leq X_n \leq m + 2\sigma$  ».

Vérifier que, pour  $n = 5$ , l'événement  $E_5$  correspond, pour un échantillon de taille 5, à un nombre de personnes possédant le gène A actif compris entre 1 et 4.

2. On désigne par  $Z_n = \frac{X_n - m}{\sigma}$ , la variable aléatoire centrée réduite correspondant à  $X_n$ .

a. Écrire, à l'aide de  $Z_n$ , l'événement  $E_n$ .

b. De quelle forme sont les valeurs  $z_k$ , pour  $k$  entier allant de 0 à  $n$ , prises par la variable aléatoire  $Z_n$  ? Quel est l'écart entre deux valeurs consécutives  $z_k$  et  $z_{k+1}$  ?

c. On souhaite représenter l'histogramme normalisé de la variable aléatoire  $Z_n$ . Il s'agit d'un histogramme pour lequel l'aire de chaque rectangle est égal à la probabilité  $P(Z_n = z_k)$ .

Quelle est la largeur commune des rectangles ?

Quelle est la hauteur du rectangle correspondant à  $Z_n = z_k$  ?

d. Sur GeoGebra, créer un curseur  $n$  allant de 5 à 5 000 puis calculer  $m$  et  $\sigma$ .

Créer l'histogramme normalisé de  $Z_n$  en entrant dans la barre de saisie :

H=Histogramme[Séquence[(i-m)/σ-0.5/σ,i,0,n+1],

Séquence[σ\*Binomiale[n,0.4,k,false],k,0,n]]

Qu'observe-t-on lorsque  $n$  augmente ?

e. Pour  $n = 5$ , à quelle aire correspond la probabilité  $P(E_5)$  ?

f. Sur GeoGebra, la fonction floor fournit l'entier directement inférieur à un nombre et la fonction Binomiale[n,p,k,true] calcule la probabilité qu'une variable aléatoire de loi binomiale de paramètres  $n$  et  $p$  prenne une valeur inférieure ou égale à  $k$ .

On suppose que  $m - \sigma$  n'est pas entier. Montrer que l'on calcule  $P(E_n)$  en saisissant :  
 $P = \text{Binomiale}[n, 0.4, \text{floor}(m + 2 * \sigma), \text{true}] - \text{Binomiale}[n, 0.4, \text{floor}(m - \sigma), \text{true}]$

Que vaut  $P(E_5)$  ?

**B. Courbe de Gauss et loi normale centrée réduite**

1. La courbe de Gauss, correspondant à la densité normale centrée réduite, représente

la fonction  $f$  définie pour tout nombre réel  $x$  par  $f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ .

Représenter  $f$  sur le fichier GeoGebra. Que constate-t-on ?

2. Calculer, à l'aide de GeoGebra, l'intégrale  $I = \int_{-1}^2 f(x) dx$ . Interpréter graphiquement cette intégrale.

3. Comparer  $P(E_n)$  et  $I$  lorsque  $n$  augmente.

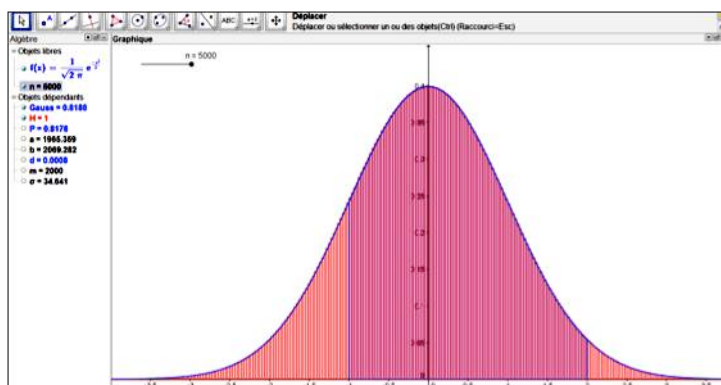
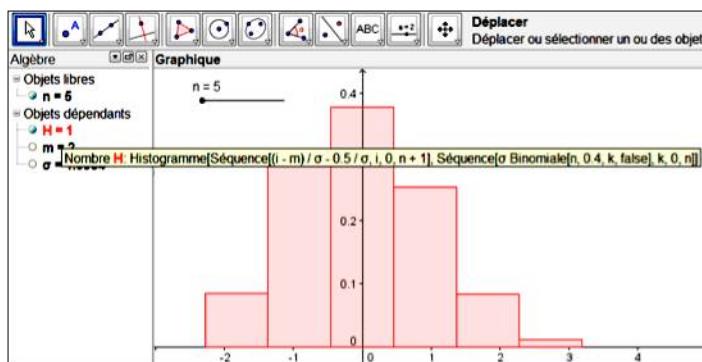


Figure 7 – Introduction du théorème de Moivre-Laplace

*Expérimentation de la notion d'intervalle de confiance*

La simulation est un outil précieux favorisant la compréhension de la notion d'intervalle de confiance et en particulier le sens du terme « confiance ».

Le fichier tableur simule des sondages aléatoires de taille 1 000 le jour de l'élection de Barack Obama c'est-à-dire avec une fréquence d'opinions favorables égale à  $p = 0,55$  dans la population.

Pour chaque sondage simulé, fournissant une fréquence d'opinions favorables  $f$ , est représenté « l'intervalle de confiance » au niveau de confiance de 95% :

$$\left[ f - \frac{1}{\sqrt{1000}} ; f + \frac{1}{\sqrt{1000}} \right].$$

- La fréquence obtenue par Barack Obama le jour de l'élection est-elle toujours comprise dans l'intervalle de confiance ?
- En faisant plusieurs simulations, estimer approximativement le pourcentage d'intervalles de confiance contenant le résultat de l'élection.
- Deux intervalles de confiance peuvent-ils être disjoints ?

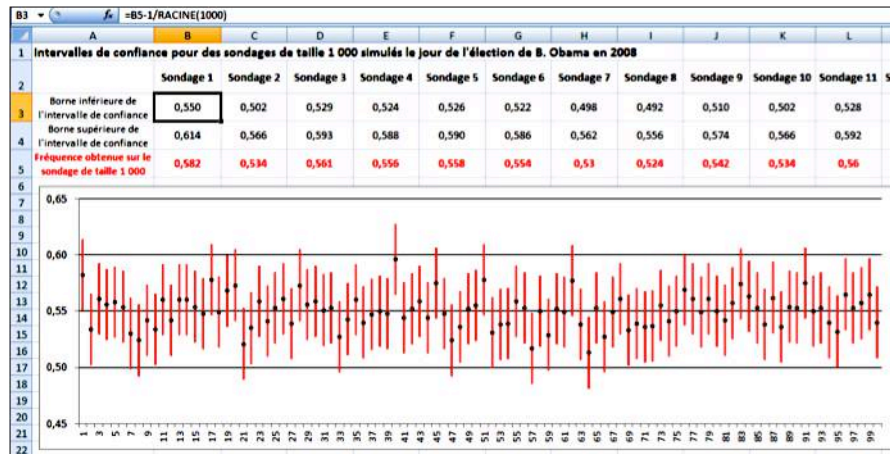


Figure 8 – Expérimentation de la notion d'intervalle de confiance

### Erreur de décision lors de la comparaison de deux intervalles de confiance

On trouve dans plusieurs manuels de terminale S (2012), dans le cas de deux intervalles de confiance disjoints, une expression du type « on conclut au niveau de confiance de 95% que les proportions correspondantes sont différentes ». Cette expression n'est pas correcte. Il suffit d'expérimenter pour s'en rendre compte.

On considère deux médicaments A et B dont on veut comparer les effets. À partir de deux échantillons de  $n$  malades, les uns traités avec le médicament A, les autres traités avec le médicament B, on construit deux intervalles de confiance à 95% des probabilités de guérison de chaque médicament. On considèrera qu'il existe une différence significative entre les deux traitements lorsque les deux intervalles de confiance sont disjoints.

Prenons par exemple  $p = 0,7$  et estimons la probabilité d'obtenir deux intervalles de confiance à 95% disjoints sur deux échantillons de taille 100 prélevés dans la même population où  $p = 0,7$ .

1. On considère les fréquences  $f_1$  et  $f_2$  de malades guéris observées sur chacun des échantillons de taille  $n$ . Montrer que les deux intervalles de confiance à 95% correspondants (on utilise l'expression au programme de terminale) sont disjoints lorsque  $|f_1 - f_2| > 0,2$ .

2. À l'aide de simulations, estimer la probabilité de l'erreur de décision correspondant à l'obtention de deux intervalles de confiance à 95% disjoints alors que la proportion dans la population est la même,  $p = 0,7$ .



3. On reprend la même étude en considérant l'intervalle de confiance à 95% « standard » utilisé dans le post-bac et correspondant à l'expression  $f$



On a obtenu les images d'écran de la feuille de calcul suivantes.



Figure 9 – Images d'écran de la feuille de calcul

- Quelle formule peut-on entrer en cellule B104 ?
- Quel est le rôle de la formule  $=SI(MAX(B104;C104)>MIN(B105;C105);1;0)$  entrée en cellule B106 ?
- Sur un grand nombre de simulations, la cellule C107 affiche en moyenne environ 0,5%. Que peut-on en déduire ?

### 2.3. En évaluant l'ensemble des compétences acquises

Prenons l'exemple du programme de terminale scientifique (2012) dont les objectifs, en termes de compétences, sont les suivants :

Outre l'apport de connaissances nouvelles, le programme vise le développement des compétences suivantes :

- mettre en œuvre une recherche de façon autonome ;
- mener des raisonnements ;



- avoir une attitude critique vis-à-vis des résultats obtenus ;
- communiquer à l’écrit et à l’oral.

Les modes d’évaluation prennent des *formes variées, en phase avec les objectifs poursuivis. En particulier, l’aptitude à mobiliser l’outil informatique dans le cadre de la résolution de problèmes est à évaluer.*

*À l’écrit au baccalauréat : algorithmique Pondichéry S 2012*

La forme écrite de l’évaluation au baccalauréat limite le cadre de l’évaluation, notamment concernant la mobilisation de l’outil informatique pour la résolution de problème. Prenons l’exemple de la question d’algorithmique posée au baccalauréat S 2012 à Pondichéry.

Un groupe de 50 coureurs, portant des dossards numérotés de 1 à 50, participe à une course cycliste qui comprend 10 étapes, et au cours de laquelle aucun abandon n’est constaté.

À la fin de chaque étape, un groupe de 5 coureurs est choisi au hasard pour subir un contrôle antidopage. Ces désignations de 5 coureurs à l’issue de chacune des étapes sont indépendantes. Un même coureur peut donc être contrôlé à l’issue de plusieurs étapes.

1. À l’issue de chaque étape, combien peut-on former de groupes différents de 5 coureurs ?
2. On considère l’algorithme ci-dessous dans lequel :
  - « rand(1, 50) » permet d’obtenir un nombre entier aléatoire appartenant à l’intervalle [1 ; 50]
  - l’écriture «  $x := y$  » désigne l’affectation d’une valeur  $y$  à une variable  $x$ .

Variables	$a, b, c, d, e$ sont des variables du type entier
Initialisation	$a := 0 ; b := 0 ; c := 0 ; d := 0 ; e := 0$
Traitement	Tant que $(a = b)$ ou $(a = c)$ ou $(a = d)$ ou $(a = e)$ ou $(b = c)$ ou $(b = d)$ ou $(b = e)$ ou $(c = d)$ ou $(c = e)$ ou $(d = e)$ Début du tant que $a := \text{rand}(1, 50) ; b := \text{rand}(1, 50) ; c := \text{rand}(1, 50) ;$ $d := \text{rand}(1, 50) ; e := \text{rand}(1, 50)$ Fin du tant que
Sortie	Afficher $a, b, c, d, e$
3. À l’issue d’une étape, on choisit au hasard un coureur parmi les 50 participants. Établir que la probabilité pour qu’il subisse le contrôle prévu pour cette étape est égale à 0,1.
4. On note  $X$  la variable aléatoire qui comptabilise le nombre de contrôles subis par un coureur sur l’ensemble des 10 étapes de la course.
  - (a) Parmi les ensembles de nombres suivants, lesquels ont pu être obtenus avec cet algorithme :  
 $L_1 = \{2, 11, 44, 2, 15\} ; L_2 = \{8, 17, 41, 34, 6\} ;$   
 $L_3 = \{12, 17, 23, 17, 50\} ; L_4 = \{45, 19, 43, 21, 18\} ?$
  - (b) Que permet de réaliser cet algorithme concernant la course cycliste ?

- (a) Quelle est la loi de probabilité de la variable aléatoire  $X$  ? Préciser ses paramètres.
- (b) On choisit au hasard un coureur à l’arrivée de la course. Calculer, sous forme décimale arrondie au dix-millième, les probabilités des événements suivants :
  - il a été contrôlé 5 fois exactement ;
  - il n’a pas été contrôlé ;
  - il a été contrôlé au moins une fois.

Figure 10 – Question d’algorithmique, baccalauréat S, Pondichéry, 2012

*Dans le cadre de travaux pratiques informatiques : simulation d'une loi normale par la méthode du rejet*

L'aptitude à mobiliser l'outil informatique dans le cadre de la résolution de problèmes doit être évaluée, en cours d'année, dans le cadre de séances de travaux pratiques. Un compte-rendu de TP peut être relevé en fin de séance, ou rédigé hors la classe. Le TP peut comprendre, à différents moments clés, des appels au professeur pour évaluer, « en direct » et à l'oral, des compétences de prise d'initiative ou d'argumentation.

On cherche à simuler des réalisations d'une variable aléatoire  $Z$  de loi normale centrée réduite, de densité  $f$  définie, pour tout réel  $x$ , par  $f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$ . La

fonction  $f$  est représentée ci-dessous.

Dans ce TP, on suppose que le générateur de nombres aléatoires de l'ordinateur simule parfaitement une variable aléatoire de loi uniforme sur l'intervalle  $[0, 1]$ .

1. a. Donner la valeur  $m$  du maximum de  $f(x)$ .
- b. Quelle est la probabilité que  $Z$  prenne ses valeurs en dehors de l'intervalle  $[-5, 5]$  ?
- c. Que vaut, approximativement, l'aire comprise entre l'axe des abscisses, la courbe représentative de  $f$  et les droites d'équation  $x = -5$  et  $x = 5$ , exprimée en unités d'aires ?

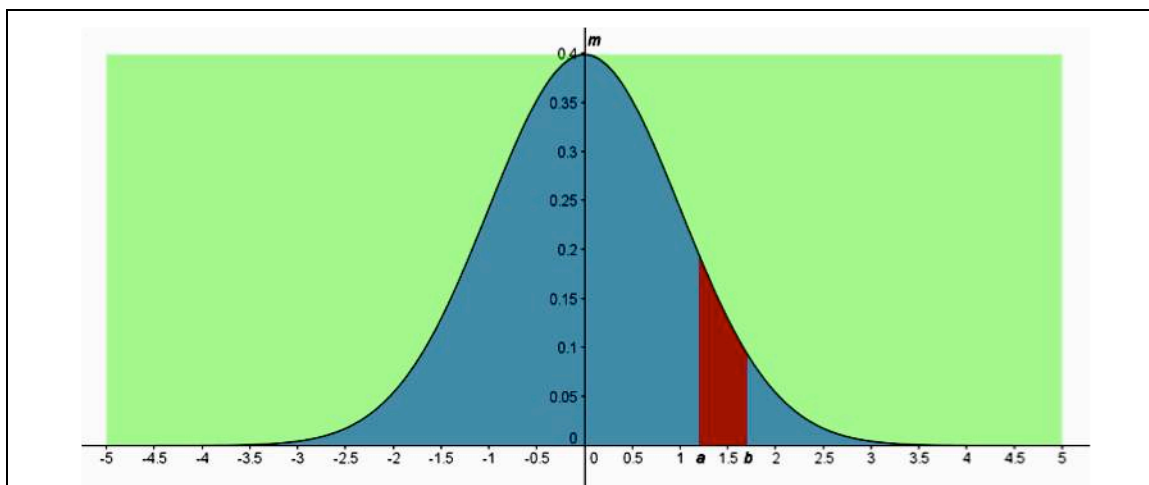


Figure 11 – Que vaut approximativement l'aire du domaine en rouge ?

2. On considère l'algorithme ci-dessous.

```

1 m=1/sqrt(2*pi)
2 function y=f(x)
3   y=m*exp(-x*x/2)
4 endfunction
5 x=-5+10*rand()
6 y=m*rand()
7 while y>f(x)
8   x=-5+10*rand()
9   y=m*rand()
10 end
11 disp(x)
    
```

- Quelle est la loi de la variable aléatoire  $X$  dont  $x$ , dans l'algorithme ci-dessus, est une réalisation ?
  - Quelle est la loi de la variable aléatoire  $Y$  dont  $y$ , dans l'algorithme ci-dessus, est une réalisation ?
  - Donner une interprétation graphique d'une réalisation  $(x, y)$  des deux variables aléatoires précédentes.
  - Interpréter graphiquement la condition de « rejet » figurant dans la boucle « tant que ».
  - Quelle est la probabilité de « rejet », c'est-à-dire que la condition de la boucle soit satisfaite ?
  - Soit  $a$  et  $b$  deux nombres de l'intervalle  $[-5, 5]$  avec  $a \leq b$ . Quelle est la probabilité qu'une valeur  $x$  acceptée (c'est-à-dire sortant de la boucle « tant que ») soit comprise dans l'intervalle  $[a, b]$  ?  
Que peut-on en déduire ?
- Implanter l'algorithme sur un ordinateur.
  - Modifier l'algorithme de sorte qu'il génère 10 000 valeurs. Implanter ce nouvel algorithme et, selon les possibilités du logiciel, afficher un histogramme et comparer avec la représentation graphique de  $f$ .

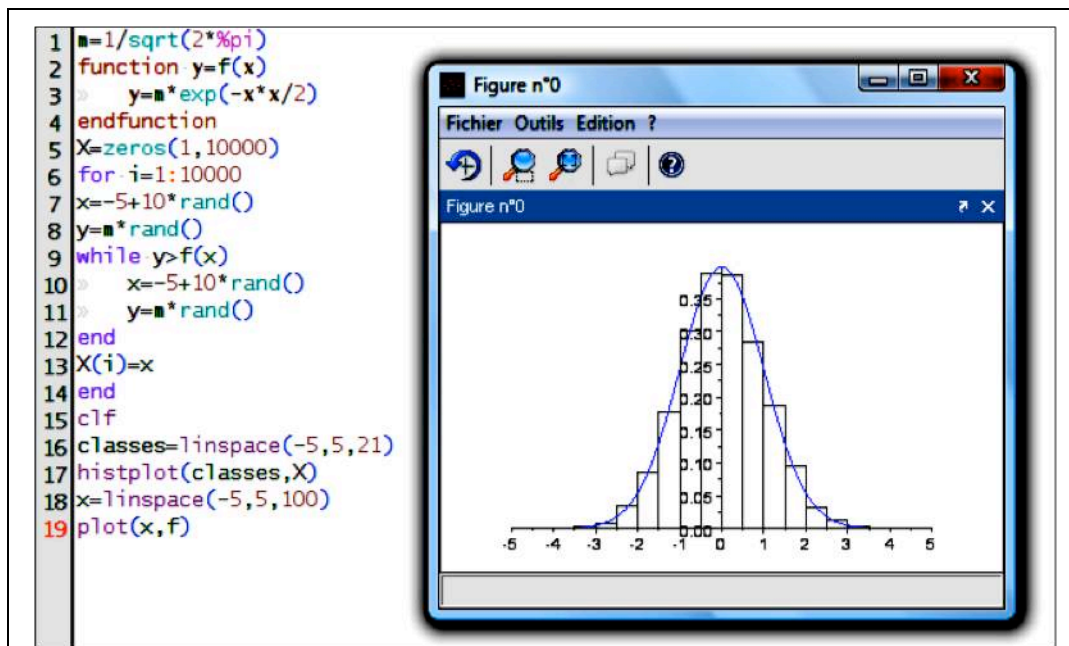


Figure 12 – L'algorithme génère 10 000 valeurs, représentées par un histogramme

Une autre méthode de simulation justifiant l'appellation de loi « normale » (illustration du théorème central limite) peut faire l'objet d'un TP analogue à l'exemple précédent :

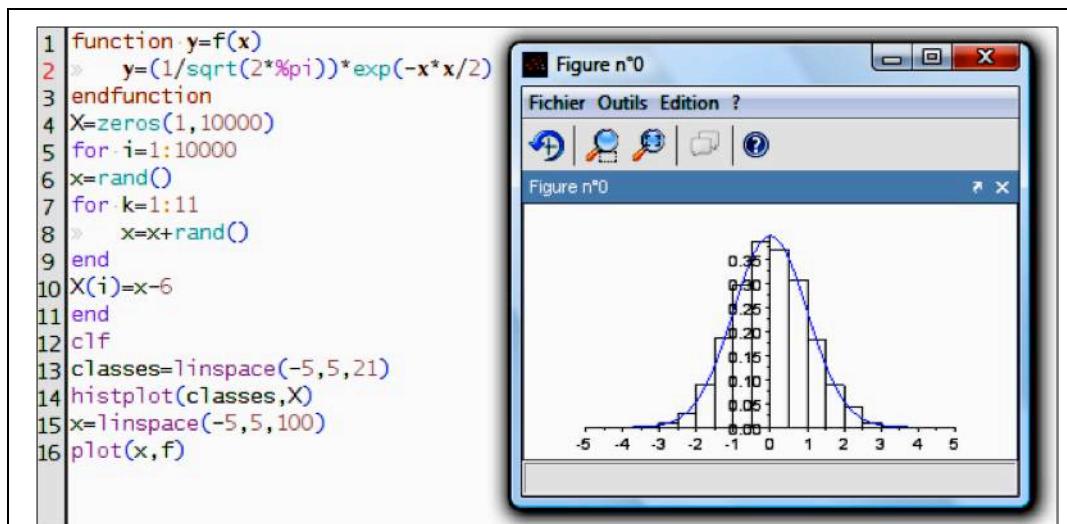


Figure 13 – Une autre méthode de simulation justifiant l'appellation de loi « normale »

Dans le cadre de la résolution d'une « tâche complexe » : les méfaits du tabac

Voici un exemple de problème présenté de façon « ouverte ».

Dans un numéro datant de 1950 du *Journal de l'Association médicale américaine*, Ernst L. Wynder, étudiant en médecine, et Ewerts Graham, chirurgien, comparent 605 cas de cancers du poumon chez les hommes à 780 hommes « témoins », n'ayant pas de cancer du poumon, recrutés dans plusieurs hôpitaux des États-Unis.

La fréquence  $f_1$  d'hommes fumeurs est de 91,3% parmi les « cas » atteints d'un cancer du poumon et la fréquence  $f_2$  de fumeurs est de 65,3% parmi le groupe « témoin ».

Wynder et Graham concluent que « l'utilisation excessive et prolongée du tabac, en particulier de cigarettes, semble être un facteur important capable d'induire le cancer du poumon ».

Comment peut-on justifier cette affirmation ?

De la seconde à la terminale, la simulation ou les outils statistiques permettent de répondre de différentes façons.

### 3. Formation des enseignants

#### 3.1. Un état des lieux en statistique

Dans le cadre d'une thèse soutenue en 2005<sup>17</sup>, Floriane Mathieu-Wozniak interrogeait 41 professeurs stagiaires de mathématiques sur leur image des différents domaines enseignés. Il en ressort que la statistique est plutôt « mathématique » : là-dessus, il n'y a pas de doute – même si elle l'est un peu moins que les autres domaines considérés. En revanche, elle est sensiblement moins brillante, quelque peu terne, elle est moins belle et moins profonde.

<sup>17</sup> Voir [3].

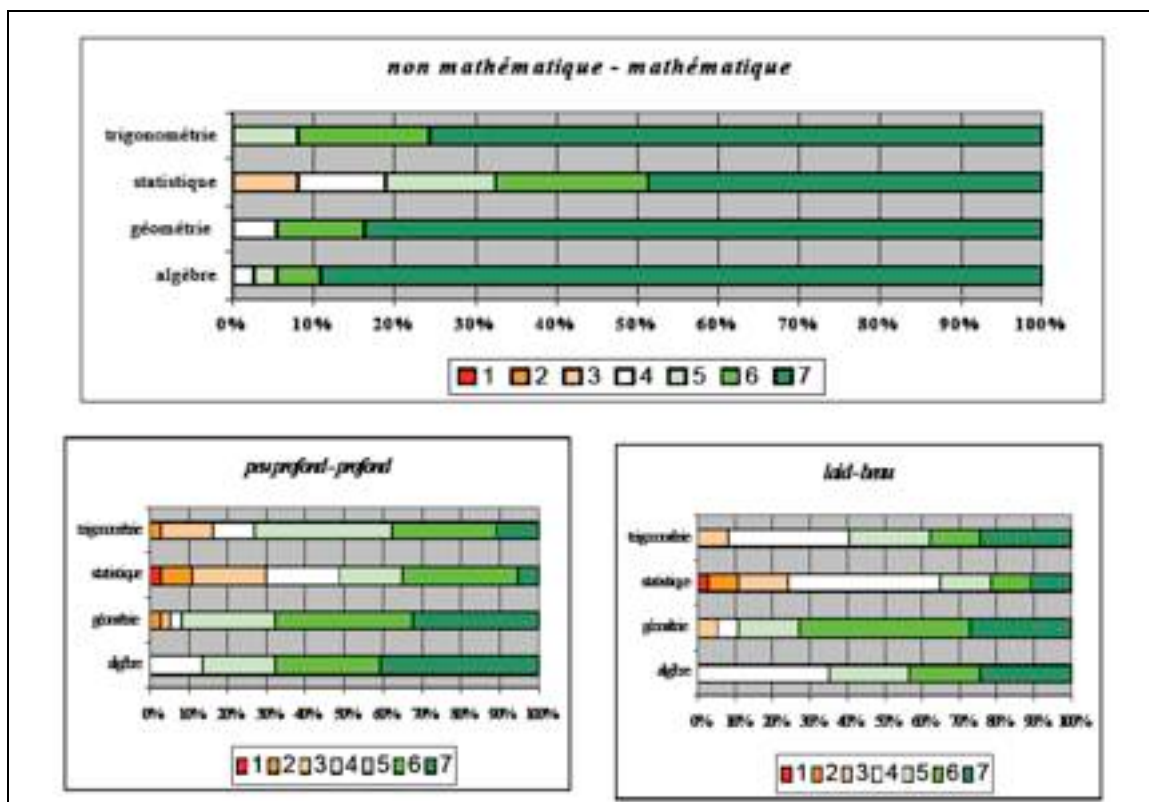


Figure 14 – Réponses de 41 professeurs stagiaires concernant leur image de différents domaines enseignés.

Cette image, assez négative, de la statistique est sans doute encore relativement majoritaire chez les professeurs de mathématiques français, même si l'on peut penser qu'elle évolue et que, par leur formation, les nouveaux professeurs l'ont en grande partie dépassée (la préparation du CAPES intègre dorénavant des leçons de statistique inférentielle issues des programmes de BTS). Elle s'explique en partie par un manque de formation dans le domaine et des contenus d'enseignement souvent techniques et peu intéressants. Par ailleurs, il s'agit en statistique de mathématiques « différentes » (rôles des raisonnements inductifs et déductifs, de la modélisation) conduisant à des attitudes pédagogiques inhabituelles.

La situation dans les classes a sans doute quelque peu évolué depuis 2005. La simulation (au programme depuis 2000) est globalement enseignée en seconde. On rencontre cependant encore fréquemment des professeurs rejetant l'enseignement de la statistique et des probabilités en fin de progression et abordant de manière incomplète les contenus du programme. L'approche fréquentiste de la notion de probabilité en troisième est peu développée et il n'est pas rare de voir, dans les cahiers des élèves, des cours de probabilités débutant par une définition correspondant au rapport des cas favorables aux cas possibles. Un enseignement trop académique de la statistique inférentielle en terminale scientifique, contextualisant peu ou mal, sans laisser de place à l'expérimentation, est à craindre, notamment de la part de professeurs peu familiers de ce type d'enseignement et non formés.



### 3.2. Des priorités de formation

Les priorités de la formation des professeurs du secondaire pour l'enseignement de la statistique pourraient être les suivantes.

1) La présentation de situations dans des contextes variés, notamment avec de « vraies données statistiques » et une initiation à la modélisation.

Il peut s'agir d'accompagner les enseignants dans les choix des activités proposés par les manuels, mais, pour l'essentiel cet apport documentaire proviendra d'institutions : formateurs IUFM, IREM, SFdS, APMEP, ressources Internet telles que celles du site Statistix... Il s'agit d'avoir les bons supports d'activité et cela demande une certaine expertise.

2) Une formation à l'organisation et à l'encadrement de la démarche d'investigation, à sa mise en œuvre avec différents logiciels et à son évaluation (prise d'information « en direct » lors de travaux pratiques, comptes-rendus de TP, possibilités offertes par les logiciels - tableurs, GeoGebra, Scilab, R ...).

On rejoint ici des besoins plus globaux de formation, qui ne sont pas spécifiques à l'enseignement de la statistique et des probabilités, mais qui, dans ce domaine, sont incontournables.

3) Des apports théoriques permettant le recul nécessaire à l'enseignement des notions nouvellement introduites dans les programmes du secondaire.

C'est une évidence, mais parce que c'est une évidence nous avons tenu à placer les deux premiers points avant. En effet, les apports théoriques ne suffisent pas et l'opérationnalité de cet enseignement passe avant tout par la contextualisation, la modélisation et l'expérimentation. Par ailleurs, la culture historique des enseignants concernant la statistique est souvent réduite. Une meilleure connaissance de la genèse, dans la première moitié du XX<sup>e</sup> siècle, des notions de test d'hypothèse et d'intervalle de confiance, des débats et controverses qu'elles ont suscités, notamment entre Fisher et Pearson-Neyman, éclaire utilement leur enseignement. La formation devrait intégrer cet aspect<sup>18</sup>.

## Conclusion

On attribue à Herbert George Wells (1866-1946) la « prophétie » suivante :

Statistical thinking will one day be as necessary for efficient citizenship as the ability to read and write.

Ce jour est venu.

L'American Statistical Association liste six recommandations pour l'enseignement de la statistique<sup>19</sup> :

1. Mettre l'accent sur les compétences statistiques du citoyen (*statistical literacy*) et développer l'esprit statistique.
2. Utiliser de vraies données.

---

<sup>18</sup> Voir [2].

<sup>19</sup> Voir [1].

3. Privilégier la compréhension des concepts plutôt que l'accumulation de connaissances techniques.
4. Favoriser l'apprentissage actif dans la classe.
5. Utiliser la technologie pour développer la compréhension des concepts et l'analyse des données.
6. Utiliser les évaluations pour améliorer l'apprentissage des élèves.

On ne peut que souscrire à ces recommandations. Les nouveaux programmes du secondaire sont relativement novateurs et ambitieux dans le domaine de la statistique, un effort particulier de formation et d'accompagnement des enseignants est nécessaire à la réussite de leur mise en œuvre.

#### REFERENCES

- [1] AMERICAN STATISTICAL ASSOCIATION, *Guidelines for Assessment and Instruction in Statistics Education – College report*, ASA 2005 et 2010. Rapport téléchargeable à l'adresse : <http://www.amstat.org/education/gaise/>
- [2] ARMATTE, Michel, *Le rôle de l'histoire dans l'enseignement de la statistique*, Revue électronique *Statistique et enseignement*, SFdS 2010, [www.statistique-et-enseignement.fr](http://www.statistique-et-enseignement.fr).
- [3] MATHIEU-WOZNIAK, Floriane, *Conditions et contraintes de l'enseignement de la statistique en classe de seconde générale. un repérage didactique*, thèse sous la direction d'Yves CHEVALLARD (2005).  
[http://tel.archives-ouvertes.fr/docs/00/06/41/60/PDF/these\\_wozniak\\_floriane.pdf](http://tel.archives-ouvertes.fr/docs/00/06/41/60/PDF/these_wozniak_floriane.pdf)
- [4] MEUSNIER, Norbert, *Sur l'histoire de l'enseignement des probabilités et des statistiques dans Histoire de probabilités et de statistiques*, Ellipse, 2004.
- [5] Ressources pour la classe de première générale et technologique – *Statistiques et probabilités*, DGESCO 2011.
- [6] Ressources pour la classe terminale générale et technologique – *Probabilités et statistique*, DGESCO 2012.

Les documents ressources sont accessibles par :

<http://eduscol.education.fr/cid45766/mathematiques-pour-le-college-et-le-lycee.html>