Cours 6 : Tableaux de contingences et tests du χ^2

Table des matières

Section 1.	Attribut moyen vs. répartition d'attributs	. 2
Section 2.	Test non paramétrique sur les fréquences	. 2
2.1.	Structure du test	. 3
Section 3.	Tableaux de contingences et fréquences marginales	. 5
9 Écrire u	ıne interprétation des données	. 6
Section 4.	Test sur la variance	. 9
© Postula	ts exacts vs. postulats asymptotiques	11
Section 5.	Conclusion	12
Exer	cices	13

Lectures

Lecture suggérée : Howell, chapitre 6, 6.1 à 6.7 inclusivement.

Objectifs

Pouvoir réaliser un test sur des répartitions dans une liste et dans un tableau avec des marges. Pouvoir réaliser un test sur la variance.

Section 1. Attribut moyen vs. répartition d'attributs

Pour faire des statistiques, il faut avant tout mesurer un attribut sur des individus. Souvent par la suite, on rapporte la mesure de tendance centrale: la moyenne. Or certains attributs ne se moyennent pas. Par exemple, lorsque l'on dit qu'il y a 30% de fumeurs, on ne dit pas que l'individu moyen fume à 30%. L'individu typique (si on entend par typique, qui représente la majorité) ne fume pas. De même, lorsqu'on rapporte que 50% des fumeurs meurent du cancer du poumon, ça ne veut pas dire que le fumeur typique est à moitié mort (quoique...). Ces deux statistiques sont basées sur des mesures de type nominales (échelle de type I). Comme on l'a vu au cours 1, une variable se mesure à l'aide d'une échelle soit nominale, ordinale, relative ou absolue. Chacune de ces échelles permet des opérations mathématiques différentes. En psychologie, plus particulièrement dans la recherche descriptive, une grande partie des données (ex. le sexe, l'origine ethnique, etc.) provient des variables classificatoires pouvant être mesurées uniquement par l'échelle nominale.

Lorsque les données sont mesurées sur des échelles de type I, la seule option possible est de classifier les données (combien sont de tel sexe? combien sont albanais? etc.) puis de compter le nombre de représentants dans chaque catégorie. Il est fréquent aussi de rapporter ces nombres en pourcentage. Cependant, dans les tests qui vont suivre, il faut utiliser la valeur des effectifs, pas les pourcentages.

Section 2. Test non paramétrique sur les fréquences

L'échelle nominale ne permet pas de tests statistiques aussi élaborés qu'une moyenne, mais elle permet quand même une comparaison entre les effectifs. La comparaison de ces décomptes permet de vérifier si les catégories sont représentées dans des proportions égales, ou dans des proportions attendues.

Par exemple, à l'automne 2000, le cours de statistique était composé de 123 femmes et de 35 hommes. Puisqu'il s'agit d'un cours obligatoire, la répartition de cet échantillon constitue peut-être un bon indice de la distribution des sexes dans le choix de cette profession. Une question naturelle alors est de savoir si la psychologie attire également les étudiants des deux sexes. Puisque notre échantillon est composé d'un total de 158 personnes, l'hypothèse nulle dicte alors une répartition attendue de 79 – 79 dans chaque catégorie.

On peut alors dresser un tableau des fréquences observées (les effectifs) que l'on note O_i dans chacune des catégories i qui nous intéressent (ici, deux catégories, hommes et femmes). En parallèle, on a aussi les valeurs attendues, notées a_i . On note O en gras car il s'agit d'une variable aléatoire qui peut changer d'un échantillon à l'autre, et a sans gras car il s'agit d'une valeur prédite à priori par notre hypothèse. Concrètement, on obtient :

	Hommes	Femmes	total
observé O	35	123	158
attendu a	79	79	158

Un autre exemple dans lequel le nombre de classe est supérieur à deux. Un journal rapporte que durant la semaine précédente, il s'est produit 427 accidents de la route, répartis ainsi pour chaque jour de la semaine :

	dim	lun	mar	mer	jeu	ven	sam	total
О	31	60	46	64	57	82	87	427
а	61	61	61	61	61	61	61	427

Encore ici, l'hypothèse que l'on veux examiner est une répartition égale des accidents au cours des jours de la semaine. Donc, a_i vaut 427 / 7 = 61.

2.1. Structure du test

Une approche est de tester si les écarts entre les valeurs observées \mathbf{O} et attendues a sont significatives. Si nous utilisons cette approche et faisons la somme des écarts, nous avons cependant un gros problème, car nous avons vu (au cours 2) que la somme des écarts donne toujours zéro. Une approche alternative est d'utiliser la somme des écarts mis au carré. De cette façon, tous les écarts deviennent positifs. Pour bien faire, on peut aussi pondérer l'écart observé par la valeur attendue a_i . On obtient une formule générale pour évaluer l'écart aux valeurs attendues :

$$\sum_{i} \frac{(\mathbf{O}_{i} - a_{i})^{2}}{a_{i}}$$

Maintenant, il est démontré que \mathbf{O} est normalement distribué quand l'effectif observé est raisonnablement grand. En effet, si a_i est, disons, 5, on s'attend à observer aussi souvent 4 que 6 (symétrie), par pur hasard. De plus, il est démontré que la variance dans le nombre d'effectifs d'une classe dont le résultat réel est a_i est aussi de a_i (c'est à dire qu'on s'attend en moyenne à une différence de $\pm \sqrt{a_i}$). Autrement dit, la somme ci-haut est une somme de scores z, des scores normalisés. Nous avons vu au cours 3 quelle est la distribution théorique d'une somme de scores normalisés, la χ^2 . Nous avons tous les ingrédients pour construire un test statistique.

a.1. Postulats

Le test est basé sur le postulat que les fréquences observées O_i sont normalement distribué et que chaque fréquence attendue est suffisamment grande. Ici, on entend par grande un $a_i > 5$. Si vous prédisez des classes avec moins de cinq effectifs, vous devez alors les regrouper pour que l'effectif prédit de cette super classe dépasse 5.

a.2. Hypothèses et seuil

Dans notre premier exemple, l'hypothèse nulle prédit une répartition égale des hommes et des femmes en psychologie. Nous pouvons alors écrire :

$$H_0: \mathbf{O}_{\text{homme}} = \mathbf{O}_{\text{femmes}}$$

 $H_1: \mathbf{O}_{\text{homme}} \neq \mathbf{O}_{\text{femmes}}$

Dans le cas des accidents automobiles, l'idée est la même, mais la formulation formelle plus longue :

$$H_0$$
: $O_{dim} = 61$ et $O_{lun} = 61$ et ... $O_{sam} = 61$
 H_1 : $O_{dim} \neq 61$ ou $O_{lun} \neq 61$ ou ... $O_{sam} \neq 61$

Notons qu'un test des fréquences est nécessairement bidirectionnel car il ne peut pas prévoir un effectif uniquement plus petit (ou uniquement plus grand).

Nous adoptons dans ces exemples un seuil usuel de 5%.

*a.*3. *Chercher le test*

Le test de fréquence est de la forme :

Rejet de H₀ si
$$\sum_{i} \frac{(\mathbf{O}_{i} - a_{i})^{2}}{a_{i}} > s(\alpha)$$

où la valeur $\sum_{i} \frac{(\mathbf{O}_{i} - a_{i})^{2}}{a_{i}}$ est distribuée comme un χ^{2} (nombre de classe – 1). Dans notre

exemple sur le sexe des futurs psychologues, le nombre de classes est 2, et donc la valeur critique, après inspection dans la table, est $s(\alpha)$ = 3.841. Pour l'exemple des accidents de la route, il y a 7 classes, et la valeur critique est $s(\alpha)$ = 12.592.

a.4. Appliquer le test et conclure

Dans le premier exemple, nous trouvons :

$$\sum_{i} \frac{(\mathbf{O}_{i} - a_{i})^{2}}{a_{i}} = \frac{(\mathbf{O}_{\text{hom }me} - a_{\text{hom }me})^{2}}{a_{\text{hom }me}} + \frac{(\mathbf{O}_{\text{femme}} - a_{\text{femme}})^{2}}{a_{\text{femme}}}$$
$$= \frac{(35 - 79)^{2}}{79} + \frac{(123 - 79)^{2}}{79} = \frac{44^{2}}{79} + \frac{44^{2}}{79} = 49.01$$

ce qui confirme ce que tous le monde sait déjà: la répartition des gars et des filles n'est pas égale en psychologie.

Dans le second exemple, on obtient, après un calcul similaire, la valeur 37.18. Comme cette valeur est supérieure à la valeur critique 12.592, on rejette H_0 et conclut que la répartition des accidents n'est significativement pas équivalente d'un jour à l'autre, et ce, avec un seuil de confiance de 5%. Puisqu'il existe au moins une différence, on peut mentionner la plus grande différence dans le tableau comme étant significative: il y a significativement plus d'accidents les samedis que les dimanches.

Précisons que dans les exemples ci-haut, nous prédisions une répartition égale. Ce n'est pas forcément notre hypothèse nulle. Une hypothèse qui prédit des valeurs attendues a_i est tout aussi possible, et est testée de la même façon que ci-haut.

Section 3. Tableaux de contingences et fréquences marginales

Supposons que dans l'étude précédente sur la répartition des sexes en psychologie, nous ayons aussi mesuré l'attitude face à la psychologie (soit favorable ou défavorable). Nous obtenons de façon générale que 63 personnes sur l'échantillon de 158 sont favorables à l'usage des statistiques en psychologie, soit un pourcentage de près de 40%. On peut se poser la question à savoir si cela est aussi vrai des hommes que des femmes. Heureusement, lors de la collecte des données sur cette deuxième question, nous avons aussi l'information sur le sexe, ce qui permet d'établir le tableau suivant, qu'on appelle un tableau de contingence :

	Femme	Homme	Total
Favorable	30 (49)	33 (14)	63
Défavorable	93 (74)	2 (21)	95
Total	123	35	158

Des marges du tableau, nous obtenons l'information précédente que les femmes sont nettement majoritaires dans le département de psychologie. De l'autre marge, nous obtenons l'information que près de 40% sont favorables à la psychologie. Est-ce qu'il en va de même pour les hommes, pour les femmes prises séparément.

L'hypothèse nulle, qui dit que 40% des femmes sont favorables à la statistique prédit donc que 0.40×123 donne le total de femmes favorables, soit 49. Inversement, 0.60×123 donne 74 femmes défavorables. Chez les hommes, 0.40×35 donne 14 alors que 0.60×35 donne 21 hommes défavorables. Nous avons mis ces valeurs entre parenthèses dans le tableau ci-haut. Comme on le voit, il semble exister des déviations importantes (surtout chez les hommes défavorables). Peut-on tester formellement cette intuition? Le test du χ^2 s'utilise aussi avec plusieurs variables classificatoires.

a.1. Postulats

Le test est basé sur le postulat que les fréquences observées O_{ij} sont normalement distribuées et que chaque fréquence attendue est suffisamment grande. Ici, on entend par grande un $a_{ij} > 5$. Si vous prédisez des classes avec moins de cinq effectifs, vous devez alors les regrouper pour que l'effectif prédit de cette super-classe dépasse 5.

a.2. Hypothèses et seuil

Dans notre exemple, l'hypothèse nulle prédit une répartition égale des attitudes des hommes et des femmes en psychologie. Nous pouvons alors écrire :

$$H_0: O_{attitude \mid homme} = O_{attitude \mid femmes} = O_{attitude}$$
 $H_1: O_{attitude \mid homme} \neq O_{attitude \mid femmes}$

où la barre verticale | se lit « étant donnée ». Autrement dit, le sexe n'influence pas l'attitude de la personne face à la statistique. Lors du cours sur les ANOVA, nous parlerons d'une "absence d'interaction" entre le sexe et l'attitude. Une façon plus courte de noter les hypothèses est:

H₀: Pas d'interaction entre le sexe et l'attitude

H₁: Interaction entre le sexe et l'attitude

Nous adoptons dans cet exemple un seuil usuel de 5%.

*a.*3. *Chercher le test*

Le test de fréquence est de la forme :

Rejet de H₀ si
$$\sum_{j} \sum_{i} \frac{(\mathbf{O}_{ij} - a_{ii})^2}{a_{ij}} > s(\alpha)$$

où la valeur $\sum_{j}\sum_{i}\frac{(\mathbf{O}_{ij}-a_{ii})^{2}}{a_{ij}}$ est distribuée comme un χ^{2} (nombre de classe i-1) × (nombre

de classe j-1). Nous avons dû utiliser une double somme car il faut faire la somme pour les deux lignes (attitudes favorable et défavorable) et pour les deux colonnes (hommes et femmes). Dans notre exemple sur le sexe et l'attitude des futurs psychologues, le nombre de colonne est 2 et le nombre de ligne, 2 aussi. On utilise donc $(2-1) \times (2-1) = 1$ degré de liberté pour rechercher la valeur critique. Elle est, après inspection dans la table, $s(\alpha) = 3.841$.

a.4. Appliquer le test et conclure

Nous trouvons:

$$\sum_{j} \sum_{i} \frac{(\mathbf{O}_{ij} - a_{ii})^{2}}{a_{ij}} = \frac{(\mathbf{O}_{HF} - a_{HF})^{2}}{a_{HF}} + \frac{(\mathbf{O}_{HD} - a_{HD})^{2}}{a_{HD}} + \frac{(\mathbf{O}_{FF} - a_{FF})^{2}}{a_{FF}} + \frac{(\mathbf{O}_{FD} - a_{FD})^{2}}{a_{FD}}$$

$$= \frac{(33 - 14)^{2}}{14} + \frac{(2 - 21)^{2}}{21} + \frac{(30 - 49)^{2}}{49} + \frac{(93 - 74)^{2}}{74}$$

$$= \frac{361}{14} + \frac{361}{21} + \frac{361}{49} + \frac{361}{74} = 55.22$$

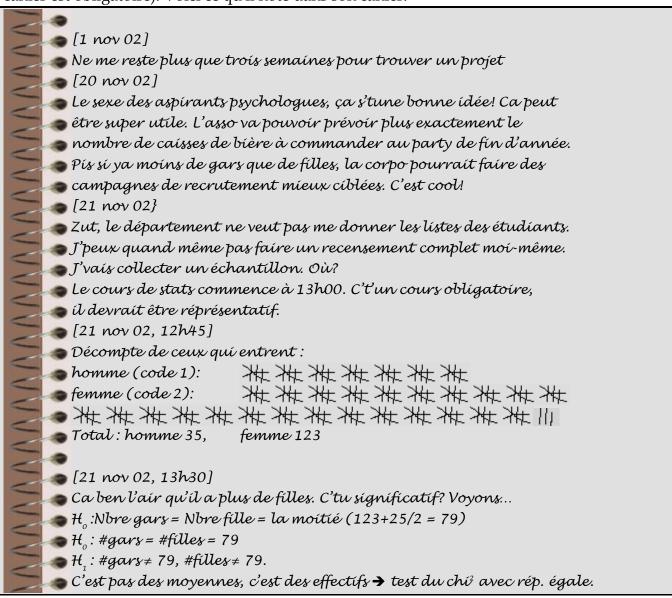
ce qui montre que les gars ont une attitude significativement différentes des filles face aux statistiques (χ^2 (1) = 55.22. p < .05). Les hommes tendent à être plus favorables aux statistiques (à plus de 90%) alors que l'inverse est vrai pour les femmes, seulement 25% d'entre elles montrant une attitude favorable aux statistiques.

9 Écrire une interprétation des données

L'écriture d'une interprétation des résultats n'est pas chose aisée. D'un côté, tout un travail de statistique a été réalisé. Or celui qui va lire votre recherche n'est pas un statisticien mais un psychologue. Vous devez expliquer les résultats en termes accessibles et significatifs pour votre lecteur. Il est probable que H_0 , μ , χ , etc. ne feront qu'égarer votre auditoire. D'un autre côté, pour des raisons de crédibilité, vous ne pouvez pas faire d'affirmations gratuites. À toutes les fois que vous rapportez une différence ou un effet, vous devez mettre dans votre rapport des signes linguistiques qui disent en substance "je n'affirme pas cela gratuitement, j'ai posé mes hypothèses et fait le test statistique approprié, et l'effet est significatif – ou pas."

Ces signes linguistiques sont les mêmes dans à peu près toutes les disciplines scientifiques: (1) l'utilisation du mot "significativement", (2) l'inclusion du résultat du test entre parenthèses, suivi du seuil α suivant cette écriture très stricte: "(nom-de-la-stat) degrés de liberté, s'il y a) = résultat, \underline{p} < seuil α)" si le test est significatif. S'il n'est pas significatif, il faut aussi rapporter la statistique, mais cette fois, " \underline{p} > seuil α)". Le signe plus petit signifie que la probabilité d'obtenir ce résultat par pur hasard est plus petit que α , ce qui veut dire qu'on a rejeté H_0 .

Voici un petit exemple qui reprend la répartition des sexes en psychologie de la section 1. Je suppose que le chercheur, un étudiant de maîtrise, tient un cahier, une sorte de journal personnel dans lequel il note ses observations et commentaires (en physique, l'usage du cahier est obligatoire). Voici ce qu'il note dans son cahier:



L'étudiant entre dans un fichier *sexe.dat*, que l'on voit dans la Figure 1 puis exécute la syntaxe suivante :

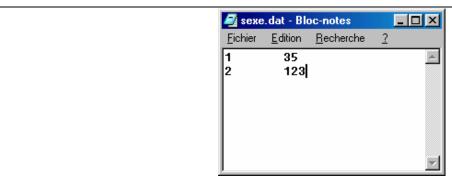


Figure 1 : Données de l'expérience sur le sexe des futurs psychologues

```
o data list file="(:\windows\bureau\sexe.dat" list
o    /sexe nbre.
o weight by nbre.
o npar test
o    /chisqu=sexe(1,2)
o    /expected=equal.
Il obtient le listing de la Figure 2 qu'il agrafe dans son cahier:
```

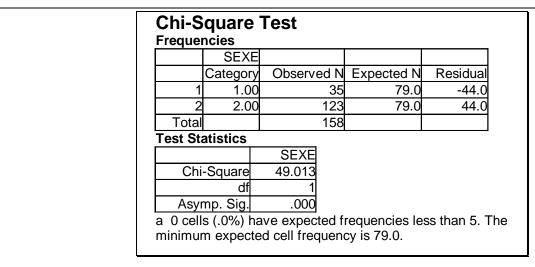


Figure 2: Listing produit pas SPSS

Voici ce qu'il écrit dans son rapport de recherche:

[Méthode]

Étant donné qu'il n'existe pas de listes d'inscrits accessible, nous avons procédé à la collecte d'un échantillon représentatif d'étudiants de psychologie à l'université de Montréal. L'échantillon est basé sur l'assistance au cours PSY 1004, un cours obligatoire pour tous les étudiants de première année.

[Interprétation des résultats]

Se trouvait présent 35 hommes et 123 femmes. Il semble qu'il y ait presque 4 fois plus de femmes que d'hommes en psychologie, ce qui est significativement différent d'une répartition égale ($\chi^2(1)$ = 49.01, p < .05).

Il se peut que notre échantillon ne soit pas absolument représentatif. Premièrement, il s'agit d'une classe de première année. Peut-être que la répartition tend à s'homogénéiser au fur et à mesure que les étudiants progressent dans leurs études. De plus, l'échantillon a été pris en fin d'année scolaire. Le taux d'absentéisme est peut-être plus grand pour les hommes que pour les femmes à cette période de l'année.

Comme on peut le voir, à part l'utilisation du mot "significativement" et la présence de la parenthèse, il n'y a aucune référence au travail de statistique dans le rapport final. Ca peut donner une impression que la statistique est négligeable en psychologie. Il n'en est rien. Si l'étudiant avait écrit entre parenthèse "($\underline{t}(177) = 4.23$, $\underline{p} > .05$)", le travail aurait été sûrement refusé et l'étudiant contraint de recommencer (ou pire, exclu des études supérieures).

Section 4. Test sur la variance

Le dernier test que nous présentons dans cette section est un test de la variance. Plus précisément, il s'agit uniquement d'un test qui permet de savoir si la variance est telle qu'attendue (selon notre hypothèse). Par exemple, soit un échantillon de 40 sujets où nous observons une variance non biaisée de 1200 unités au carré. Supposons que nous ayons une théorie qui prédit que la variance devrait être de 1000. Est-ce que la prédiction faite par cette théorie est rejetée?

Dans la pratique courante de la recherche, il est plutôt rare de voir une théorie qui fait des prédictions sur la variance. Il en existe, et en général, ce sont les plus complètes (car elles font aussi des prédictions sur la moyenne), mais elles se comptent sur le bout des doigts. Néanmoins, le test développé ici est à la base de l'ANOVA, le test statistique le plus utilisé pour distinguer des moyennes. Il est donc important de saisir la provenance de ce test, ce que nous voyons ici.

Examinons un peu la variance d'un échantillon. Si nous faisons des manipulations algébriques, nous avons (dans la suite, nous ne mettrons pas l'étiquette "n - 1") :

$$\vec{\mathbf{X}}^2 = \frac{1}{n-1} \sum_{i} (\mathbf{X}_i - \overline{\mathbf{X}})^2$$

Or,

$$(\mathbf{X}_{i} - \mu)^{2} = [(\mathbf{X}_{i} - \overline{\mathbf{X}}) + (\overline{\mathbf{X}} - \mu)]^{2}$$
$$= (\mathbf{X}_{i} - \overline{\mathbf{X}})^{2} + 2(\mathbf{X}_{i} - \overline{\mathbf{X}})(\overline{\mathbf{X}} - \mu) + (\overline{\mathbf{X}} - \mu)^{2}$$

Il s'ensuit que

$$\begin{split} \sum_{i} (\mathbf{X}_{i} - \mu)^{2} &= \sum_{i} (\mathbf{X}_{i} - \overline{\mathbf{X}})^{2} + 2 \sum_{i} (\mathbf{X}_{i} - \overline{\mathbf{X}}) (\overline{\mathbf{X}} - \mu) + \sum_{i} (\overline{\mathbf{X}} - \mu)^{2} \\ &= \sum_{i} (\mathbf{X}_{i} - \overline{\mathbf{X}})^{2} + 2 (\overline{\mathbf{X}} - \mu) \sum_{i} (\mathbf{X}_{i} - \overline{\mathbf{X}}) + \sum_{i} (\overline{\mathbf{X}} - \mu)^{2} \\ &= \sum_{i} (\mathbf{X}_{i} - \overline{\mathbf{X}})^{2} + 0 + n (\overline{\mathbf{X}} - \mu)^{2} \\ \sum_{i} (\mathbf{X}_{i} - \overline{\mathbf{X}})^{2} &= \sum_{i} (\mathbf{X}_{i} - \mu)^{2} - n (\overline{\mathbf{X}} - \mu)^{2} \end{split}$$

(on se rappelle que la somme des écarts à la moyenne donne toujours zéro – cours 2). Si l'on divise les deux côtés par σ_0^2 , où σ_0^2 est la variance prédite par la théorie,

$$\frac{1}{\sigma_0^2} \sum_i (\mathbf{X}_i - \overline{\mathbf{X}})^2 = \frac{1}{\sigma_0^2} \sum_i (\mathbf{X}_i - \mu)^2 - \frac{n}{\sigma_0^2} (\overline{\mathbf{X}} - \mu)^2$$
$$= \sum_i \frac{(\mathbf{X}_i - \mu)^2}{\sigma_0^2} - \frac{(\overline{\mathbf{X}} - \mu)^2}{\sigma_0^2 / n}$$

Maintenant, la partie de gauche $\frac{1}{\sigma_0^2} \sum_i (\mathbf{X}_i - \overline{\mathbf{X}})^2$ est égale à $(n-1)\frac{\overline{\mathbf{X}}^2}{\sigma_0^2}$, soit un ratio de la

valeur de la variance obtenue de notre échantillon divisée par la variance attendue multiplié par (n - 1). Si notre hypothèse est vraie, les deux variables sont égales, et le ratio donne 1. Donc, la partie de gauche devrait donner une valeur proche de (n - 1).

La partie de droite, si on regarde attentivement est a) une somme de scores bruts normalisé plus b) une moyenne d'échantillon normalisée par son erreur type. Si vous vous rappelez de votre cours 3, on voit que a) est un χ^2 avec n degrés de liberté et que b) est χ^2 avec un seul degré de liberté (il n'y a qu'un terme). Ceci signifie que la partie de gauche, contenant les éléments de notre hypothèse (variance observée et variance prédite) est distribuée comme une χ^2 avec n-1 degrés de liberté. Il ne nous en faut pas plus pour faire un test statistique. Le test peut être unidirectionnel ou bidirectionnel.

a.1. Postulats

Le test de la variance est basé sur le postulat que les données brutes X_i sont normalement distribuées. Il s'agit d'un postulat fort, car il ne porte pas sur la moyenne, mais sur chaque donnée brute (on ne peut pas faire appel au théorème central limite ici). Dans les faits cependant, le test semble assez résistant (si n >>), de telle façon que pour des données approximativement normales, les conclusions sont encore valables.

a.2. Hypothèses et seuil

Dans notre exemple, l'hypothèse nulle prédit que la variance réelle de la population est donnée par σ_0^2 , une valeur décidée a priori par notre théorie. Nous pouvons alors écrire :

$$H_0: \sigma^2 = \sigma_0^2$$

 $H_1: \sigma^2 \neq \sigma_0^2$

Comme toujours, nous choisissons dans notre exemple un α = 5%. La distribution du χ^2 n'est pas symétrique; il faut donc trouver deux bornes, une inférieure $s^-(\alpha/2)$ et une supérieure $s^+(\alpha/2)$.

a.3. Chercher le test

Le test de la variance est de la forme :

Rejet de H₀ si
$$(n-1)\frac{\ddot{\mathbf{X}}^2}{\sigma_0^2} > s^+(\alpha/2)$$
 ou $(n-1)\frac{\ddot{\mathbf{X}}^2}{\sigma_0^2} < s^-(\alpha/2)$

où la valeur $(n-1)\frac{\overrightarrow{\mathbf{X}}^2}{\sigma_0^2} \sim \chi^2(n-1)$. Les valeurs critiques sont: $s^+(\alpha/2) = 56.895$ et $s^-(\alpha/2) = 22.878$ avec 38 degrés de liberté (car 39 n'est pas tabulé). La figure qui suit montre les deux seuils.

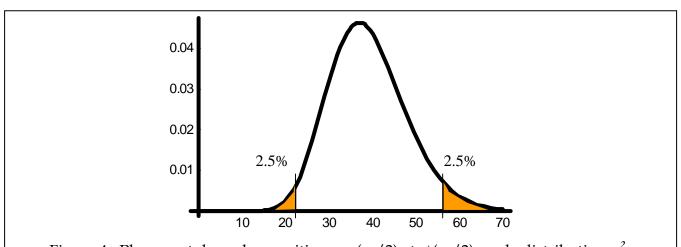


Figure 4 : Placement des valeurs critiques s- $(\alpha/2)$ et s+ $(\alpha/2)$ sur la distribution χ^2 .

a.4. Appliquer le test et conclure

Le test devient:

$$(n-1)\frac{\overrightarrow{\mathbf{X}}^2}{\sigma_0^2} = (40-1)\frac{1200}{1000} = 46.8$$

La valeur obtenue de notre échantillon n'étant pas plus grande que la valeur critique, on ne rejette pas H_0 ; la variance observée est compatible avec la prédiction du modèle.

• Postulats exacts vs. postulats asymptotiques

Certaines distributions sont basées sur le postulat que n >> alors que d'autres sont basées sur une taille n donnée. On appelle le premier type de postulat un postulat asymptotique. La distribution normale, telle que définie dans la théorie des erreurs de Gauss, est basée sur un nombre de facteur n non spécifié mais grand. Il s'agit donc d'une théorie basée sur un postulat asymptotique. La distribution de Weibull est aussi basée sur un postulat asymptotique puisque le nombre de compétiteurs n est non spécifié mais doit être

grand. À l'opposé, la distribution binomiale est basée sur un n précis (il faut connaître le paramètre n pour lire dans la table B(n, p)); il s'agit d'un postulat exact.

Les distributions asymptotiques sont plus difficiles à démontrer (le théorème central limite pris 100 ans à démontrer; le cas général pour la Weibull a été prouvé en 2002 par moimême). Cependant, ils sont beaucoup plus simples d'utilisation puisqu'on n'a pas à connaître la valeur de n. De plus, les postulats asymptotiques peuvent être utilisés pour rejeter des théories de la pensée. En effet, comme on ne connaît pas le nombre de neurones dans le cerveau (ou très approximativement), une théorie de l'esprit ne peut pas être bâtie sur un nombre n de neurones (ou de synapses, selon ce qui sera le plus pertinent).

Voici deux exemples de théorie du fonctionnement cognitif qui font appelle à des arguments asymptotiques: (1) Les *réseaux connexionnistes* classiques fonctionnent en simulant un grand nombre de "neurones" qui peuvent, après des stimulations, accroître ou réduire leurs connections réciproques. Les premiers neurones peuvent être connectés à une caméra et les derniers à un bras robot. Les connexions de ces neurones simplifiés sont donc des "facteurs" affectant la réponse, ils sont en grand nombre, et à peu près autant excitateur qu'inhibiteur. Or, suivant la théorie des erreurs de Gauss, il faut que les réponses de ce système soit distribuées de façon normal (argument asymptotique utilisant le théorème central limite). Cependant, quand on enregistre les temps de réponse des humains, ils ne sont jamais symétriques autour du temps de réponse moyen, ce qui invalide le modèle normal, et donc, les réseaux connexionnistes classiques (rejet de H₀).

Les réseaux de courses postulent plutôt que les neurones vont réagir aussitôt qu'un récepteur sensoriel ou un neurone antécédent s'active. Il s'agit donc d'un modèle de course (le plus rapide envoie son signal en premier). En utilisant un argument asymptotique, le modèle prédit est la distribution de Weibull qui prédit que les temps de réponse sont asymétriques, tout comme ceux qu'on observe. Le modèle de course n'est donc pas rejeté (non rejet de H_0).

Section 5. Conclusion

Exercices

- 1. Un chercheur fait une étude auprès de deux petits groupes d'étudiants (littérature et optométrie). Il s'intéresse à la variable « Qualité de vie » qu'il mesure suivant le revenu annuel. Lequel, parmi les tests suivant, doit-il réaliser?
 - a) Test z
 - b) Test sur la médiane
 - c) Test t
 - d) Test χ^2
- 2. En une heure, les nombres d'utilisateurs de 6 guichets automatiques sont les suivants : 12, 32, 21, 41, 67, 37. Vérifiez que ces guichets sont uniformément visités sur la base de votre échantillon avec un seuil de décision de 0.01.
- 3. Un nouveau vaccin a été testé sur 150 enfants, dont 70 dans un groupe contrôle. Six enfants traités sont malades contre 25 dans le groupe contrôle. Choisissez un seuil et vérifiez l'efficacité du vaccin.
- 4. Un candidat à la mairie de Montréal demande un sondage. Effectué auprès de 111 hommes et 133 femmes, il révèle que 49 personnes sont contre lui alors que 171 lui sont favorables, et 24 incertaines. Parmi les femmes, ces nombres deviennent 35 contre, 80 pour. Est-ce que le candidat peut compter également sur les hommes et les femmes pour être élu?
- 5. Un journal rapporte le nombre de meurtres commis à Montréal au cours des années 1984 à 1989 : 34, 27, 41, 25, 18, 35. Vérifiez l'hypothèse d'une répartition égale avec un seuil de 5%.

- 6. Un sondage effectué auprès des 12 employés d'un bureau rapporte que 4 d'entre eux ont une opinion favorable envers leur nouveau patron, 5 sont défavorables, et 3 restent neutres. Les opinions se répartissent-elles également? (Attention à la pogne)
- 7. Au cours d'une recherche, vous avez mesuré la taille de 300 universitaires québécois. Vous voulez vérifier si les données de cette variable se distribuent normalement (seuil de 1%). Pour cela, vous devez : a) établir une distribution de fréquences à parti des tailles observées (tableau donné ci-bas), et b) dresser un tableau des valeurs attendues si la distribution est belle et bien normale. La moyenne de votre échantillon est 1.625 m avec un écart type non biaisé de 7.0 cm.

Tailles	Observées	Attendues
1.45-1.50	10	?
1.50-1.55	28	
1.55-1.60	67	
1.60-1.65	84	
1.65-1.70	65	
1.70-1.75	34	
1.75-1.80	12	
	300	300

8. Une compagnie observe que les ventes du dernier mois se répartissent comme suit dans leurs 5 commerces : 81, 84, 75, 78, 82, en K\$. Doit-elle fermer une boutique (seuil de 1%)?