

T. D. n° 2

Inférence conditionnelle pour deux échantillons

Exercice 1. La suite de l'exercice 1 du T.D. n° 2

Vous souvenez-vous quelle était la conclusion de cet exercice ? La conclusion conduisait à dire que les estimations en mètres sont légèrement plus hautes que les estimations données en pieds. Ici, nous pouvons réanalyser ces données avec une optique d'inférence conditionnelle. D'abord nous allons à nouveau convertir les mètres en pieds et stocker le vecteur des observations dans une variable y .

1. Télécharger le fichier « roomwidth » du package « HSAUR »
2. Exécuter les lignes de commande suivantes :


```
> convert <- ifelse(roomwidth$unit=="feet",1,3.28)
> feet <- roomwidth$unit=="feet"
> metre <- !feet
> y <- roomwidth$width*convert
```
3. Le test statistique réalisé est simplement un test sur la différence en moyenne. Exécuter les commandes suivantes :


```
> Tobs <- mean(y[feet])-mean(y[metre])
> Tobs
```
4. Afin d'approcher la distribution conditionnelle de la statistique de test T , nous allons calculer 9999 fois la statistique de test pour des valeurs mélangées de y . Une permutation de y est obtenue avec la fonction `sample()`. Donc, maintenant que nous avons compris l'esprit, nous allons exécuter les commandes suivantes :


```
> meandiffs <- double(9999)
> for (i in 1:length(meandiffs)) {
+   sy <- sample(y)
+   meandiffs[i] <- mean(sy[feet])-mean(sy[metre]) }
> hist(meandiffs)
> abline(v=Tobs, lty=2)
> abline(v=-Tobs, lty=2)
```
5. Le graphique 1 page 2 que vous venez d'obtenir représente la distribution de la statistique de test T sous l'hypothèse nulle d'indépendance des estimations de la largeur de l'amphithéâtre et des deux groupes. Maintenant la valeur T_{obs} de la réalisation de la statistique observée pour les données originales non mélangées peut être comparée avec la distribution de T sous l'hypothèse nulle (ce qui est représenté par l'une des barres verticales dans votre graphique). La p -valeur, i.e., la probabilité que la statistique de test T soit plus grande que 8.859 ou plus petite que -8.859 est


```
> greater <- abs(meandiffs) > abs(Tobs)
> mean(greater)
```

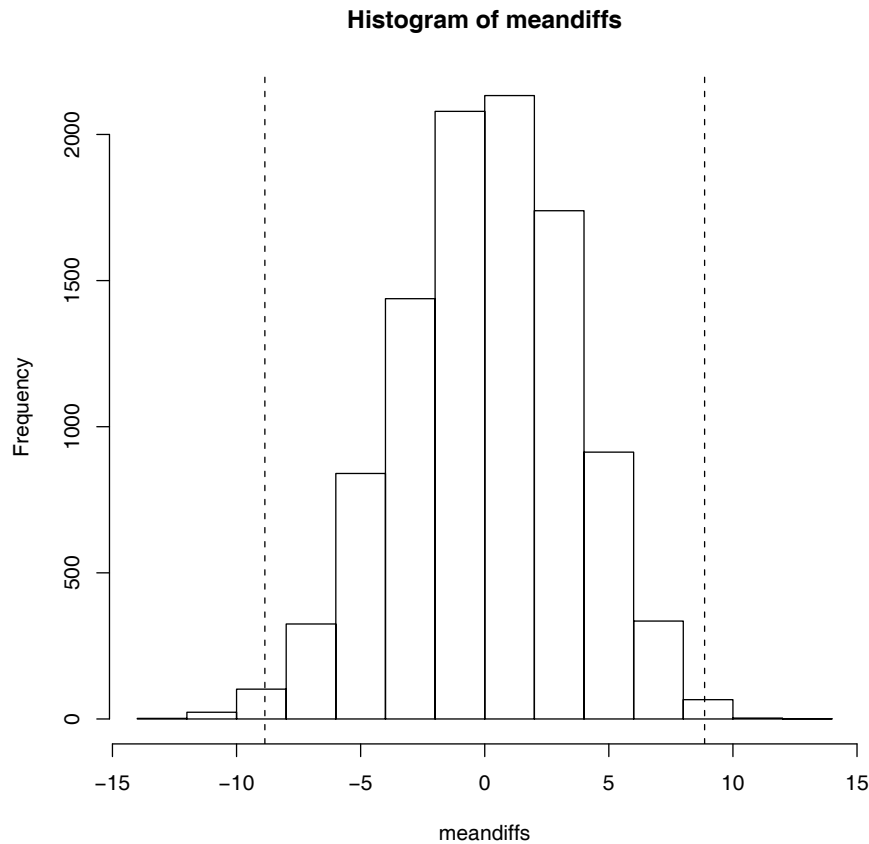


FIGURE 1. Graphique 1

6. L'intervalle de confiance est donné par la ligne de commande suivante :

```
> binom.test(sum(greater), length(greater))$conf.int
```

Remarque : Refaites ces lignes de commande depuis la question 4. Qu'observez-vous ?

7. Une autre façon de traiter le problème.

```
> library("coin")
> independence_test(y ~ unit, data=roomwidth,
+ distribution = "exact")
```

Remarque : Pour certaines situations, incluant les analyses montrées ici, il est possible de calculer la p -valeur exacte, i.e., la p -valeur basée sur la distribution évaluée sur toutes les permutations possibles des valeurs de y .

La fonction `independence_test()` (package *coin*, Hothorn et al., 2005) peut-être utilisée pour calculer la p -valeur exacte comme nous l'avons fait ci-dessus.

De façon similaire, la distribution conditionnelle exacte du test de la somme des rangs de Wilcoxon Mann-Whitney peut-être calculée par une fonction implémentée dans le package *coin*. D'ailleurs, faisons-le, c'est-à-dire exécuter la ligne de commande suivante :

```
> wilcox_test(y ~ unit, data=roomwidth, distribution="exact")
```

Que remarquez-vous ?

Exercice 2. Suicides.

Mann (1981) établit un compte rendu d'une étude qui examine les causes d'un comportement indifférent ou concerné d'une foule lorsqu'une personne menace de se suicider en sautant du haut d'un immeuble. L'hypothèse que la foule s'intéresse à la personne est plus probable en période chaude. Mann (1981) classe 21 dénombrements de tentative de suicide par deux facteurs, la période de l'année et si la foule est intéressée ou non. Les données sont présentées dans le tableau ci-dessous et la question est de savoir si ils ont raison d'invoquer cette hypothèse. Les données venant de l'hémisphère nord, la période Juin-Septembre représente la période chaude.

	Intéressée	Indifférente
Juin-Septembre	8	4
Octobre-Mai	2	7

1. Télécharger le fichier « suicides » du package « HSAUR » ou alors rentrer les données.
2. Afficher-le à l'écran.
3. Combien de lignes ? Combien de colonnes ?
4. Peut-on envisager de faire un test du χ^2 ? Les conditions sont-elles respectées ? Pour cela calculer le tableau des effectifs théoriques.
5. À la vue des résultats obtenus ci-dessus, qu'envisagez-vous de faire pour répondre à la question posée ? En fait, quel test proposez-vous de réaliser ?

Exercice 3. Tératogenèse.

Dans cet exemple, un médecin (MD) et un assistant de recherche (RA) estiment le nombre d'anomalies (0, 1, 2, 3) présentes sur un bébé pour 395 bébés. Le tableau de données se présente de la façon suivante :

	RA			
MD	0	1	2	3
0	235	41	20	2
1	23	35	11	1
2	3	8	11	3
3	0	0	1	1

1. Nous allons d'abord rentrer les données. Pour cela exécuter les lignes de commande suivantes :

```
> anomalies <- as.table(matrix(c(235,23,3,0,41,
+ 35,8,0,20,11,11,1,2,1,3,1),ncol=4,
+ dimnames=list(MD=0:3, RA =0:3)))
> anomalies
```

2. Nous allons nous intéresser à tester si le nombre d'anomalies estimées par le médecin diffère structurellement du nombre rendu par l'assistant de recherche. Comme nous comparons des observations « appariées », i.e., une paire de mesures pour chaque nouveau-né, un test d'homogénéité marginale (la généralisation du test de McNemar rencontré au T.D. numéro 5, exercice 4) peut être envisagé. Exécuter la ligne de commande suivante
`> mh_test(anomalies)`
3. Que concluez-vous à l'issue de ce test ?
4. Il est à noter que le facteur « anomalies » prend des valeurs ordonnées. Cela peut rentrer en ligne de compte, en exécutant la ligne de commande suivante :
`> mh_test(anomalies, scores=list(c(0,1,2,3)))`
5. Qu'observez-vous ? Obtenez-vous la même conclusion que précédemment ?
Remarque : Il est possible de calculer une estimation de la p -valeur exacte en utilisant l'option `distribution="approximate"` (`B=10000`).

.....