

Comparaison de « moyennes » de 2 échantillons aléatoires indépendants

Quel test utiliser ?

# Introduction

- Une question récurrente dans la comparaison de deux « moyennes » est la condition de normalité des variables analysées pour certains tests (test t)
- Plusieurs avis se confrontent:
  - 1) On ne vérifie pas la normalité et on utilise un test paramétrique
    - Quand  $n$  est « suffisamment grand », le Théorème Central Limite postule que l'estimateur de la moyenne suit une loi normale et on utilise le test t
    - Limite: qu'est-ce qu'un grand échantillon ( $n=30?$ ,  $n=50?$ )
    - Certains auteurs considèrent que les tests paramétriques sont robustes jusqu'à un certain point à la violation de la normalité
    - Limite: les violations peuvent être de différente nature (aplatissement, asymétrie)
  - 2) On ne vérifie pas la normalité et on utilise un test non paramétrique
    - Les tests non paramétriques s'appliquent quelque soit la forme de la distribution
    - Limite: les tests non paramétriques sont réputés moins puissants que les tests paramétriques si les conditions d'application sont respectées (normalité)
  - 3) On vérifie systématiquement la normalité
    - Utilisation d'un test paramétrique (test t) en cas de non rejet
    - Utilisation de tests non paramétriques ou robustes en cas d'échec
    - Limite: quel test de normalité utiliser, performance de ces tests réputée médiocre

→ L'objectif de ce travail est d'explorer les performances de ces approches

# Plan

## I. Objectifs

## II. Présentation

- des tests
- de la méthode

## III. Résultats

## IV. Conclusion / Discussion

# Objectifs (1/2)

- Explorer les erreurs de type I et II
  - de 5 différents tests :
    - Tests paramétriques : test t et test z
    - Tests non paramétriques : test de MW et test exact
    - Test robuste
  - sous diverses distributions:
    - Normale: distribution « idéale »
    - Uniforme: distribution symétrique
    - Seminormale: distribution asymétrique
    - Bimodale: le second mode correspond à des « outliers »
    - Lognormale: distribution asymétrique très skewed
  - Sous diverses tailles d'échantillon
  - Sous différentes tailles d'effet

# Objectifs (2/2)

- Comparer 2 statistiques de test: celle du t et du MW
- Explorer la puissance de différentes stratégies
  - Stratégie 1 : t test d'emblée
  - Stratégie 2 : MW d'emblée
  - Stratégie 3 :
    - Test de Normalité de Kolmogorov-Smirnov (sur chacun des 2 groupes)
    - Si rejet de l'hypothèse de normalité dans au moins l'un des deux groupes alors utilisation de MW
    - Sinon utilisation du t test

# Présentation des tests

- 1. test t
- 2. test z
- 3. test de MW
- 4. test exact
- 5. test robuste

# Présentation des tests (1/5)

## test t pour échantillons indépendants

- Condition d'application (Robert F. WOOLSON):
  - Les observations de chaque échantillon sont issues d'une distribution normale
  - Les variances des 2 échantillons sont égales
- H0: les moyennes des 2 populations normales sont égales

– Statistique de test :

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_{X_1}^2 + S_{X_2}^2}{n}}} \quad t \sim T(2n - 2)$$

# Présentation des tests (2/5)

## test Z

### – Condition d'application

- Les observations de chaque échantillon sont issues d'une distribution normale
- Les écarts types des populations source sont supposées connues

### – H0: les moyennes des 2 populations normales sont égales

### – Statistique de test

$$z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_{X_1}^2 + S_{X_2}^2}{n}}} \quad z \sim \mathcal{N}(0,1)$$



# Présentation des tests (3/5)

## test de Mann Whitney

- Pas d'hypothèse sur la forme des distributions
- Test qui se base sur le rang des observations
- H0: La probabilité qu'une observation aléatoire du 1<sup>er</sup> groupe soit plus grande qu'une observation aléatoire du 2<sup>nd</sup> groupe est 0.5 (AUC-ROC=0.5)
- Attention, H0 n'est pas l'égalité des moyennes, ni l'égalité des médianes, mais plutôt « la médiane des différences est égale à 0 »

$$K = \frac{12}{N(N+1)} \sum_{j=1}^k \left\{ \frac{R_j^2}{n_j} - 3(N+1) \right\}$$

$$K \sim \chi^2(k-1)$$

Avec:

- $R_j$ : somme totale des rangs
- $n_j$ : taille d'échantillon
- $k$ : nombre de groupes (dans notre situation,  $k=2$ )
- $N$ : taille de l'échantillon total (dans notre situation,  $N=2n$ )
- Possibilité d'ajuster la statistique de test en cas d'ex aequo

1947

On a Test of Whether one of Two Random Variables is Stochastically Larger than the Other

Author(s): H. B. Mann and D. R. Whitney

# Présentation des tests (4/5)

## test exact

- Test basé sur les permutations
  - Aucune hypothèse sur les distributions si  $n < 50$
  - Approximation normale usuelle (test Z) si  $n \geq 50$
- $H_0$ : les distributions sont identiques dans les 2 groupes

# Présentation des tests (5/5)

## test robuste

- Objectif : proposer une alternative qui soit moins impactée par les outliers ou autres violations des hypothèses du modèle (normalité)
- $H_0$ : les moyennes des 2 populations sont égales
- Principe: plus un individu est loin de la moyenne de son groupe, plus son poids sera faible (il existe plusieurs fonctions permettant de définir les poids)
- On utilise dans R la fonction `lmrob` (paramètres par défaut)
- J.W. Tukey (1979):  
“... just which robust/resistant methods you use is not important – what is important is that you use some. It is perfectly proper to use both classical and robust/resistant methods routinely, and only worry when they differ enough to matter. But when they differ, you should think hard.”

Méthode

# Méthode (1/2)

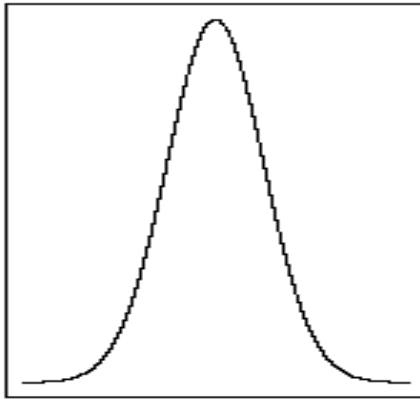
- On utilise des simulations pour répondre aux objectifs. On réplique 1000 échantillons pour chaque condition d'expérience. Les 5 tests sont alors appliqués sur chaque échantillon
- Description des simulations:
  - Génération de la population source selon une certaine distribution (N=10.000.000)
  - Par simplicité, on centre réduit ( $m=0$ ,  $sd=1$ ) les distributions générées
  - Sélection aléatoire de l'échantillon (bras 1) d'une certaine taille
  - Sélection aléatoire de l'échantillon (bras 2) de même taille auquel on rajoute une certaine différence

# Méthode (2/2)

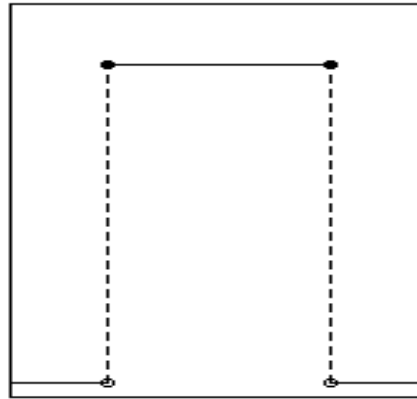
- Condition de simulations
  - Distribution : Normale, Uniforme, Seminormale, Bimodale et Lognormale
  - Taille de l'effet en Ecart Type (ET) : 0, 0.2, 0.5, 0.8 (nulle, faible, modérée, forte selon Cohen)
  - Taille d'échantillon par groupe : 10, 20, 30, 50, 100, 200, 500, 1000

# Formes des distributions étudiées

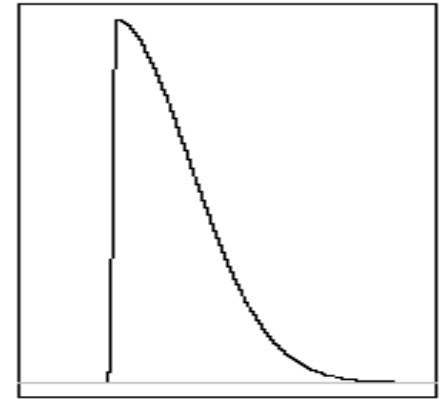
**Normal**



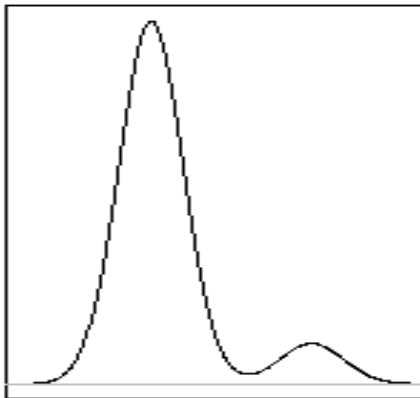
**Uniform**



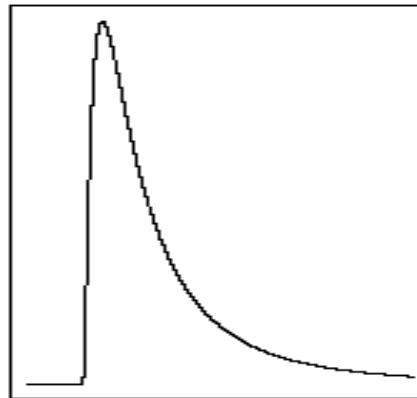
**Seminormal**



**Normal bimodal (10%~N(5;1))**



**LogNormal**

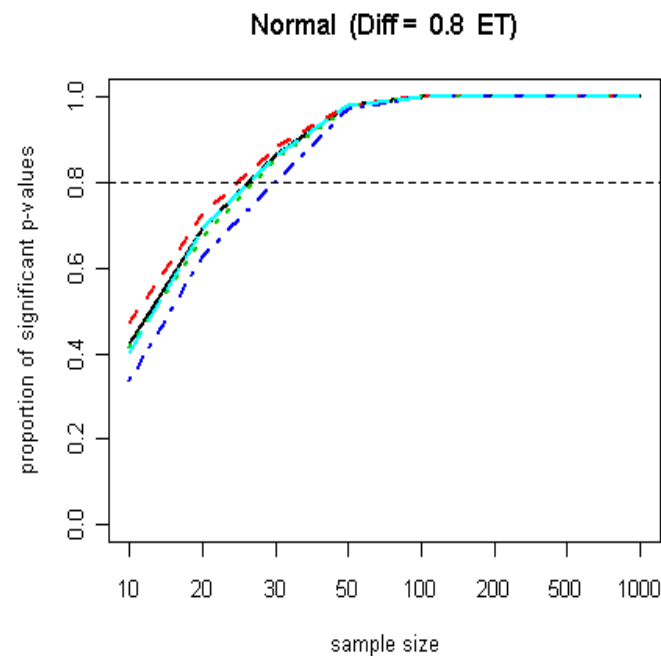
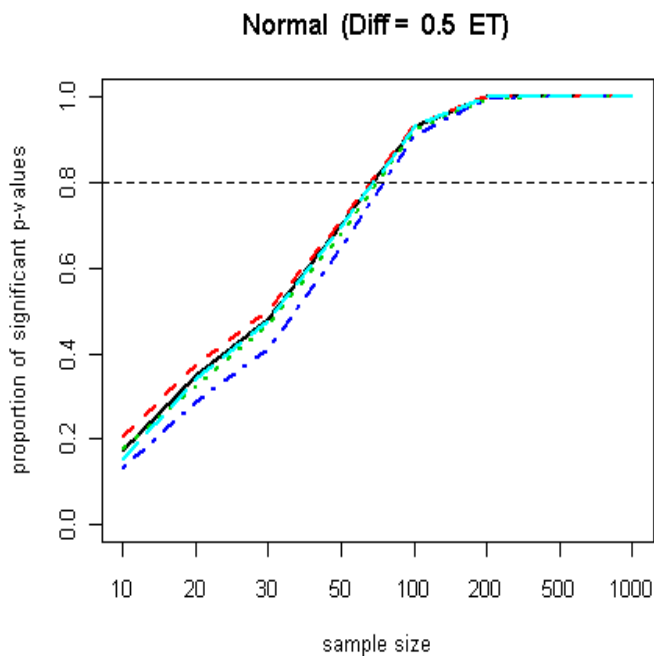
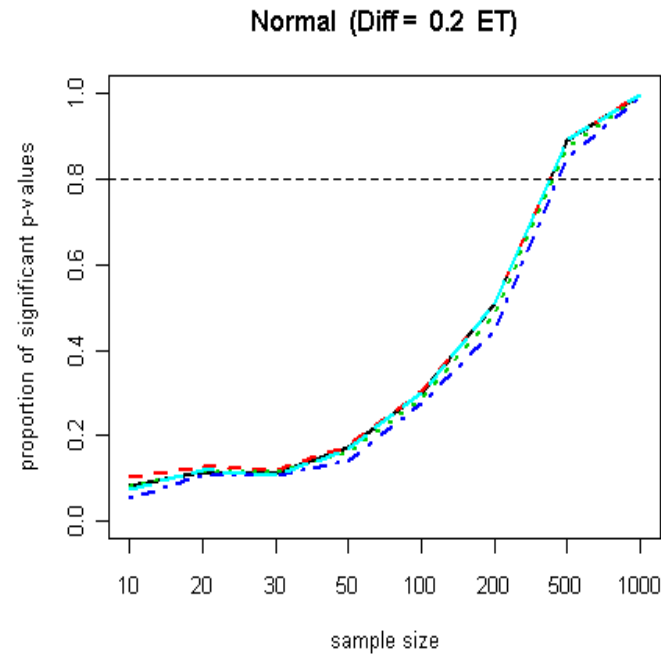
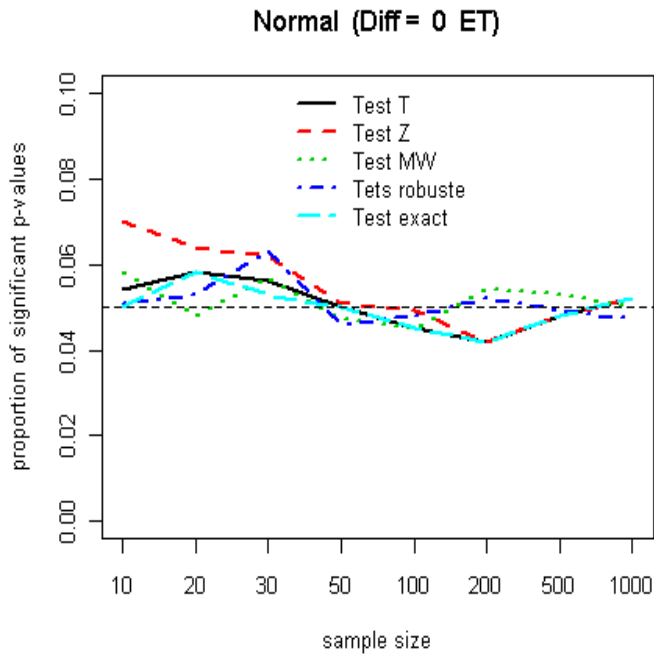


# Résultats

- Erreurs de type I et II des 5 tests, distribution par distribution



# Distribution normale

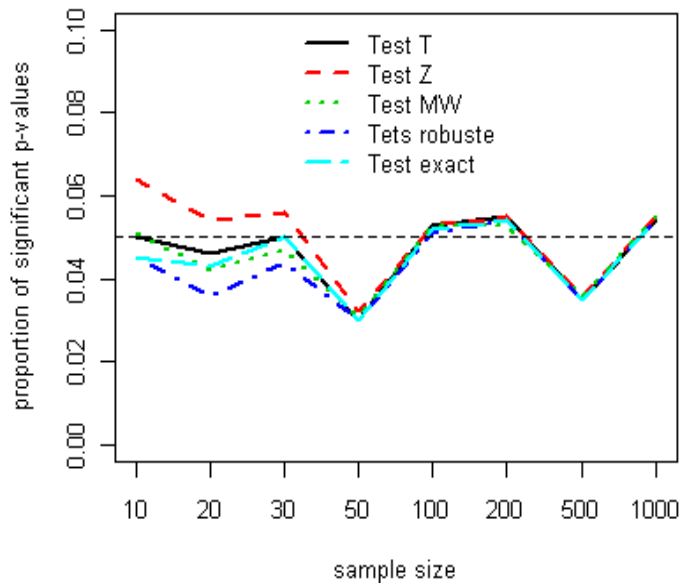


- Erreur de type 1 sensiblement équivalente

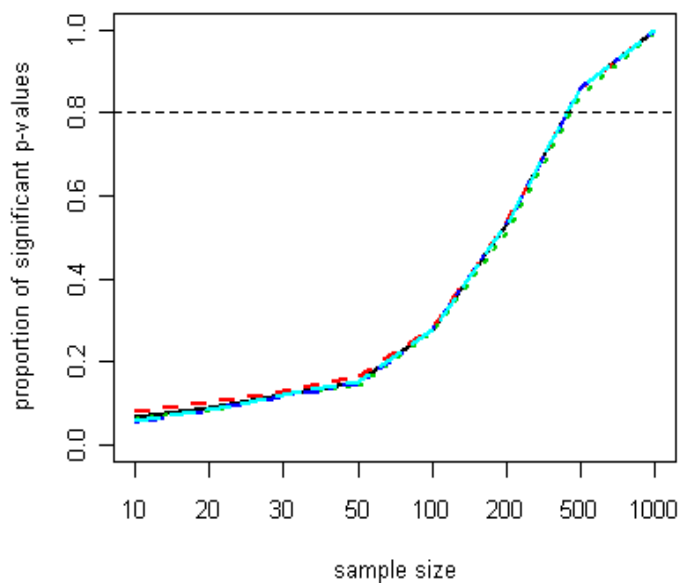
- Différence de puissance faible entre les tests

- Conclusion: les tests sont quasi équivalents pour des distributions normales

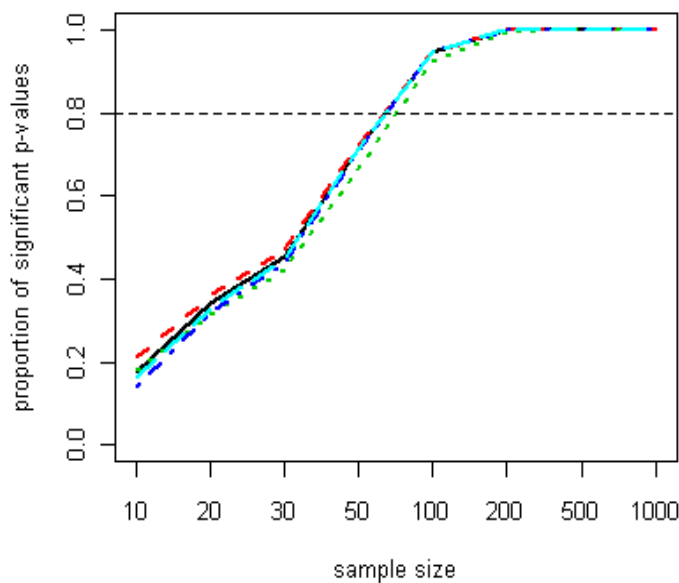
Uniforme (Diff = 0 ET)



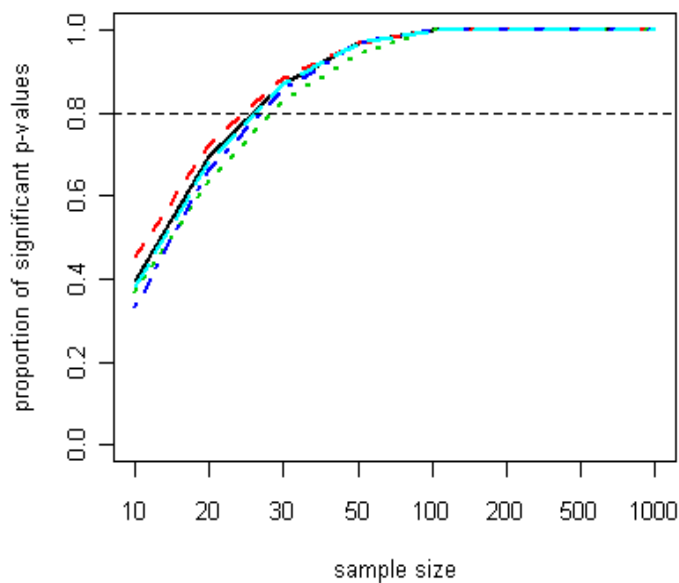
Uniforme (Diff = 0.2 ET)



Uniforme (Diff = 0.5 ET)



Uniforme (Diff = 0.8 ET)



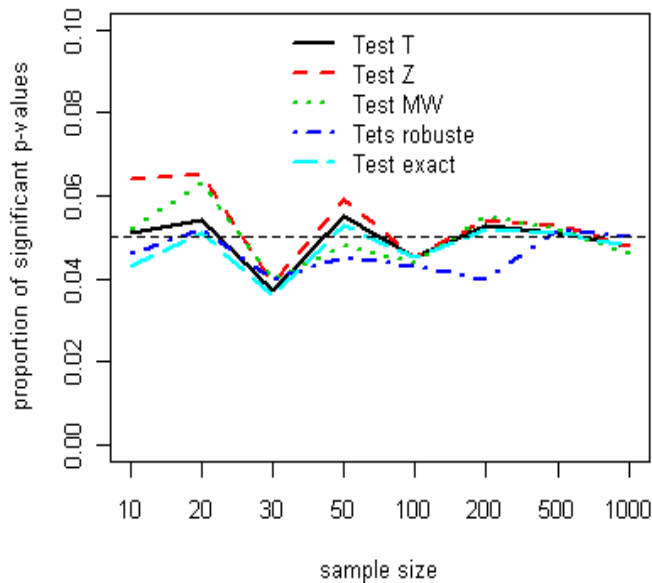
# Distribution uniforme

Alpha  $\approx$

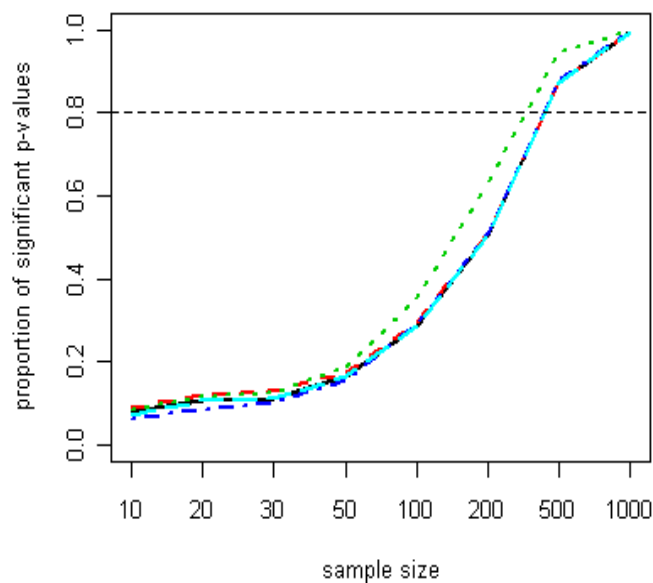
Puissance  $\approx$

les tests sont  
quasi  
équivalents  
pour des  
distributions  
uniformes

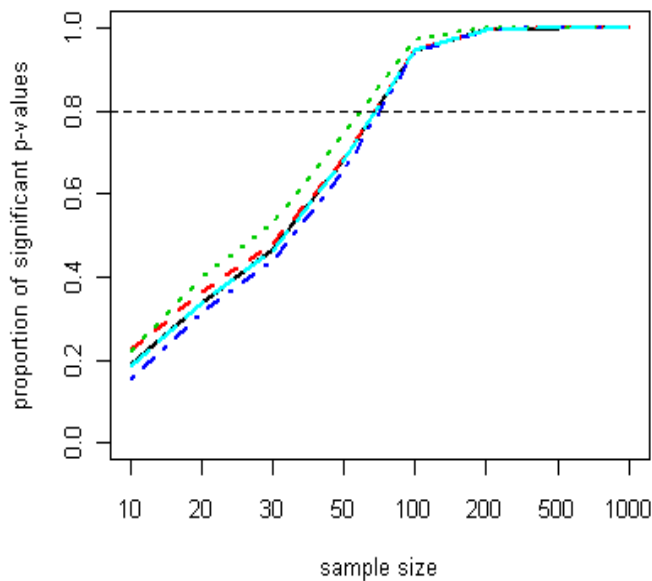
Seminormal (Diff = 0 ET)



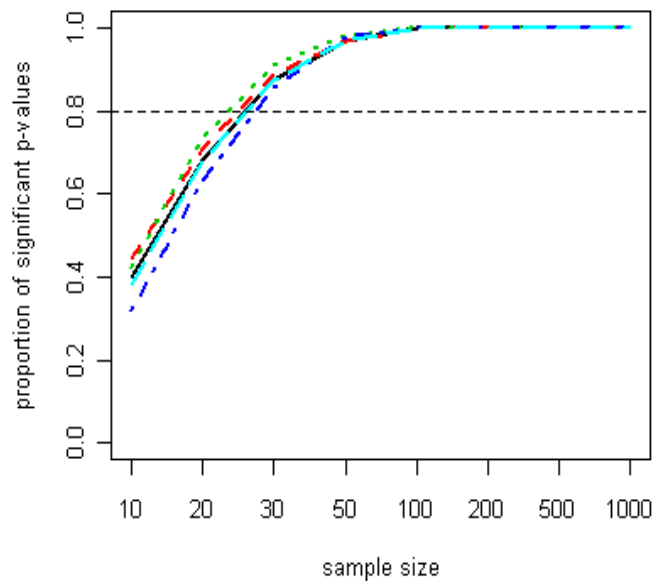
Seminormal (Diff = 0.2 ET)



Seminormal (Diff = 0.5 ET)



Seminormal (Diff = 0.8 ET)



# Distribution seminormale

Erreur de type 1 sensiblement équivalente

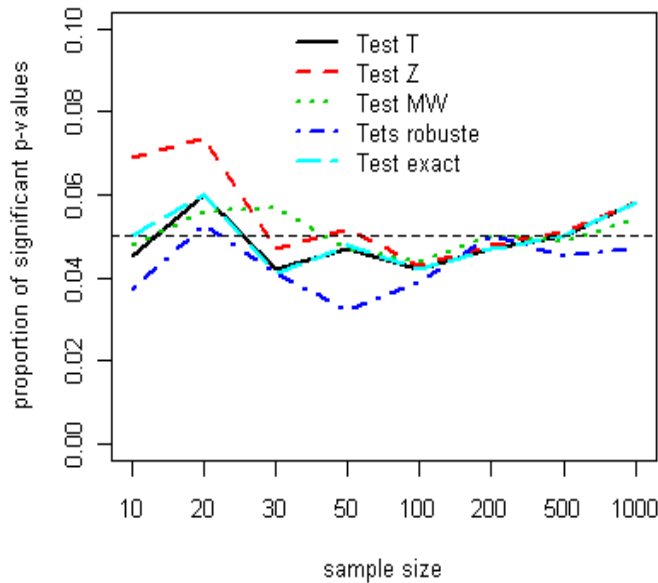
Puissance sensiblement équivalente sauf pour le test de Mann-Whitney qui est légèrement plus puissant

# Distribution normale bimodale

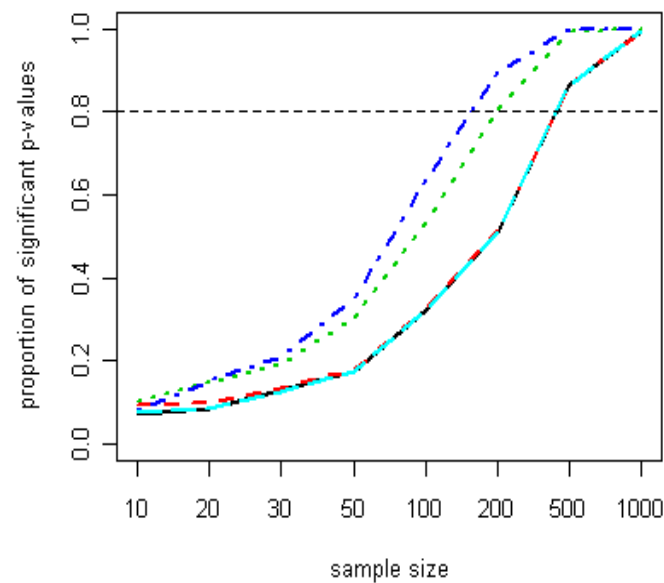
Erreur de type  
1 sensiblement  
équivalente

Test de Mann-  
Whitney et test  
robuste sont  
plus puissants

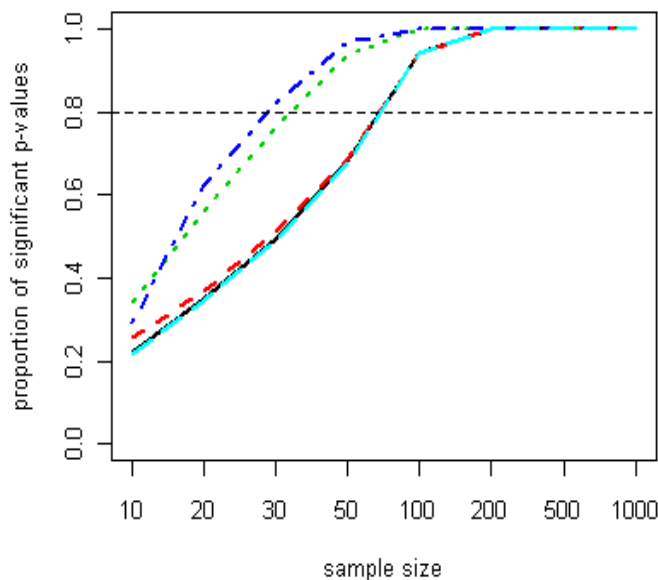
Bimodal\_5EC\_100OUT (Diff = 0 ET)



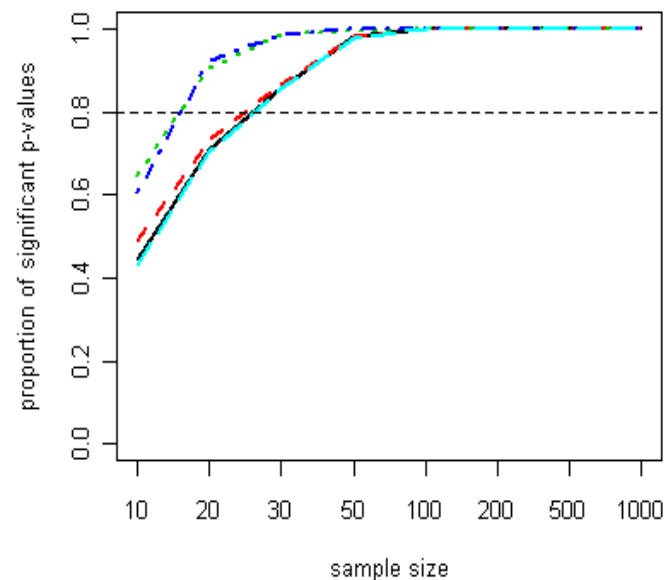
Bimodal\_5EC\_100OUT (Diff = 0.2 ET)



Bimodal\_5EC\_100OUT (Diff = 0.5 ET)



Bimodal\_5EC\_100OUT (Diff = 0.8 ET)



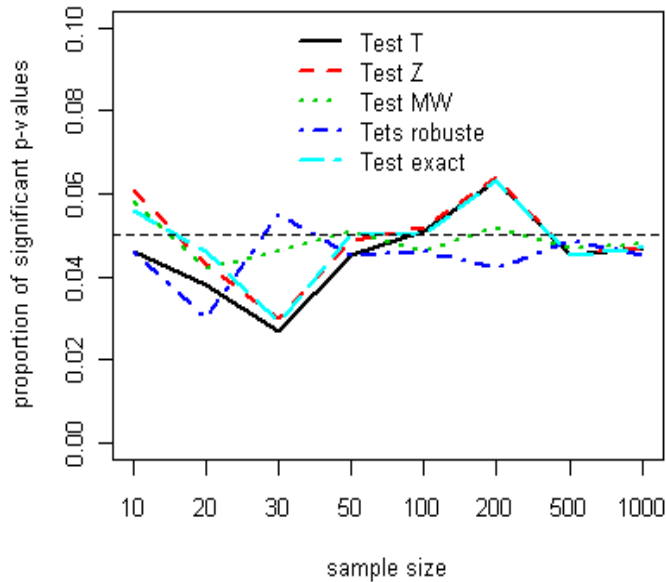
# Distribution lognormale

- Erreur de type 1 sensiblement équivalente

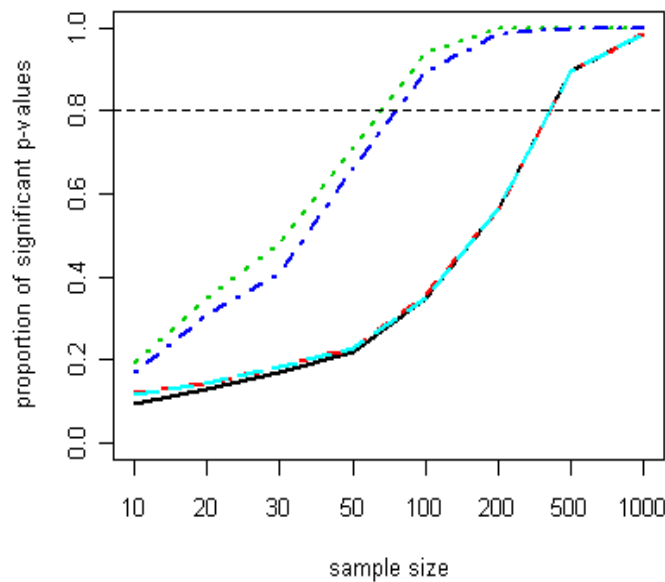
- Puissance plus forte pour le test robuste et de Mann-Whitney

- A l'inverse, le test Z est un peu + conservateur

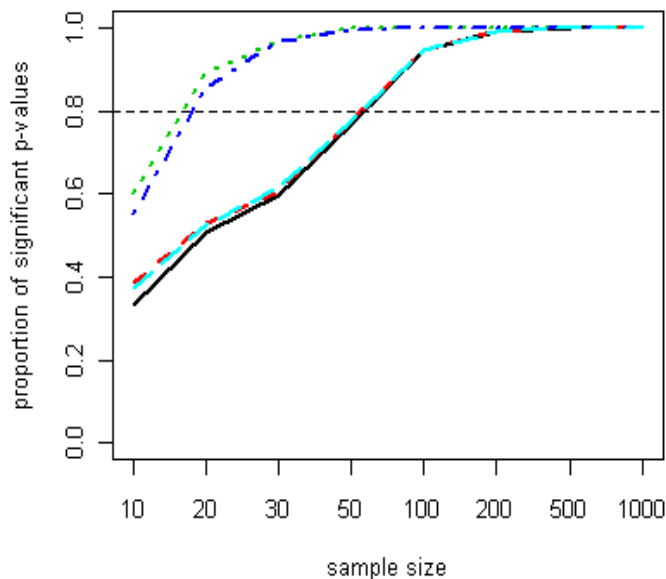
Exponential of normal (Diff = 0 ET)



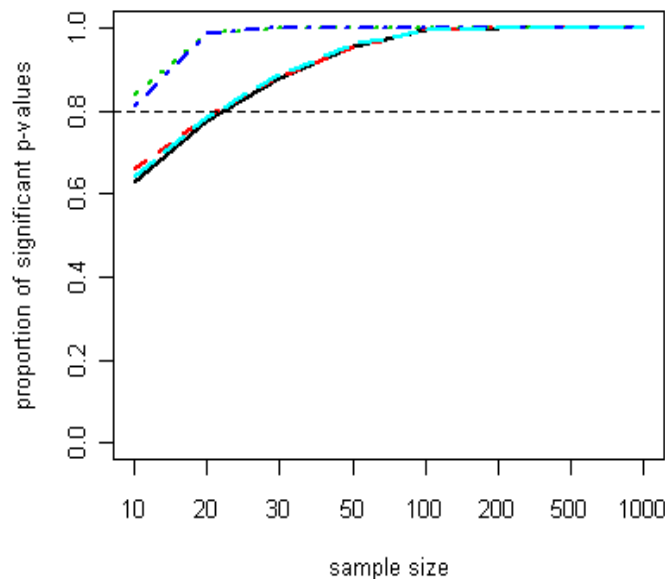
Exponential of normal (Diff = 0.2 ET)



Exponential of normal (Diff = 0.5 ET)



Exponential of normal (Diff = 0.8 ET)



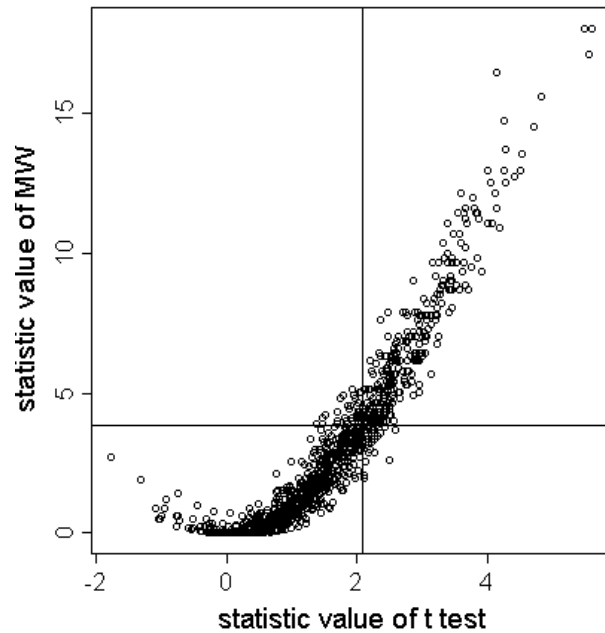
# résumé

- Sous des distributions normales: (très) faible perte de puissance du test robuste et de MW par rapport aux tests t, Z et exact
- Toute violation de la normalité n'a pas les mêmes impacts sur les erreurs de type I et de type II :
  - Distributions uniformes : résultats identiques que sous distributions normales
  - Distributions seminormale : résultats proches mais avec un faible avantage en faveur du MW
  - Distributions bimodales (outliers) : MW et robuste net gain de puissance avec avantage au robuste entre les 2
  - Distributions lognormales : net gain de puissance pour MW et robuste avec avantage au MW
- L'asymétrie et les outliers jouent sur la puissance des tests et le test de MW se comporte bien dans tous les cas de figure

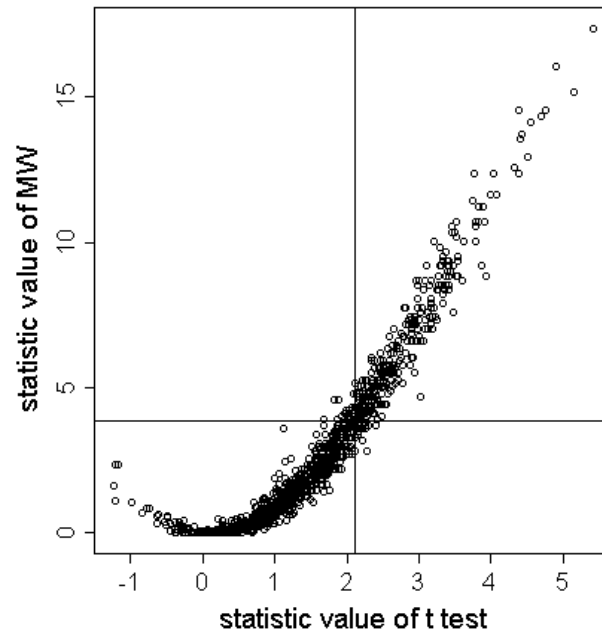
# Comparaison des statistiques de test du t test et du test de MW

- On cherche à explorer plus en détail les différences observées entre le test t et MW. On se place sous le cas  $n=20$  et  $\text{diff}=0.5$  ET
- Pour chaque simulation, on représente conjointement la statistique de test du test t et du test de MW

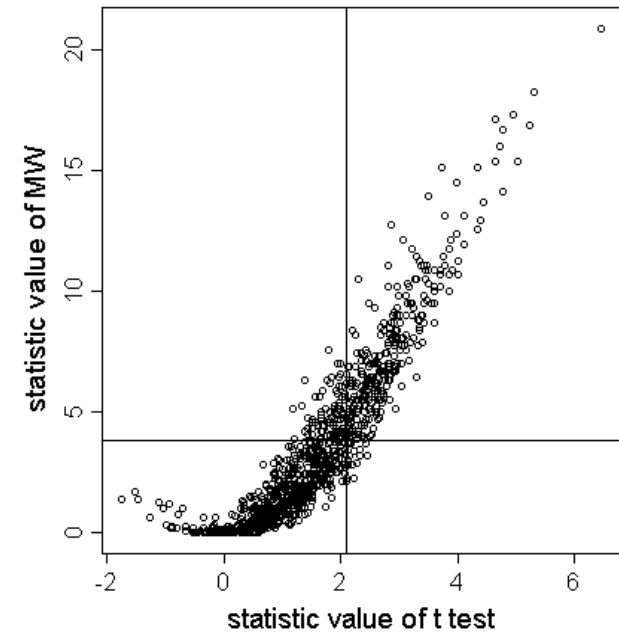
NormalD=0.5ET, N=20



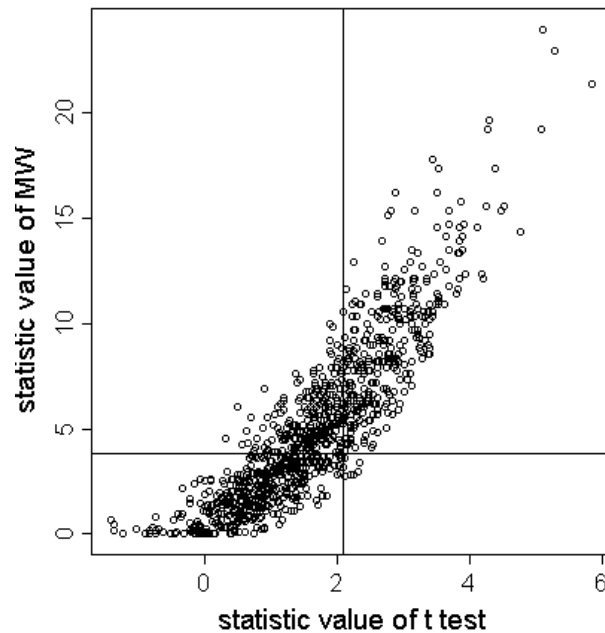
UniformedD=0.5ET, N=20



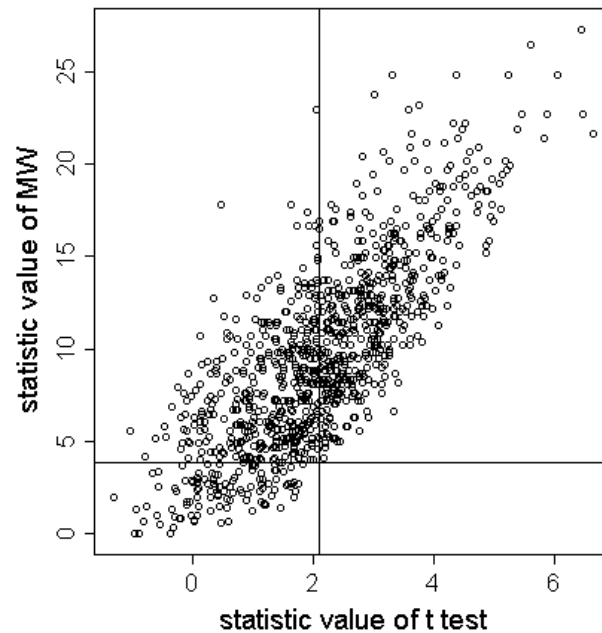
SeminormalD=0.5ET, N=20



Bimodal\_5EC\_10OUTD=0.5ET, N=20



Exponential of normalD=0.5ET, N=20



- Ligne horizontale  $h \approx 3.841$  (valeur critique MW)
- Ligne verticale  $v \approx 2.1$  (valeur critique t)



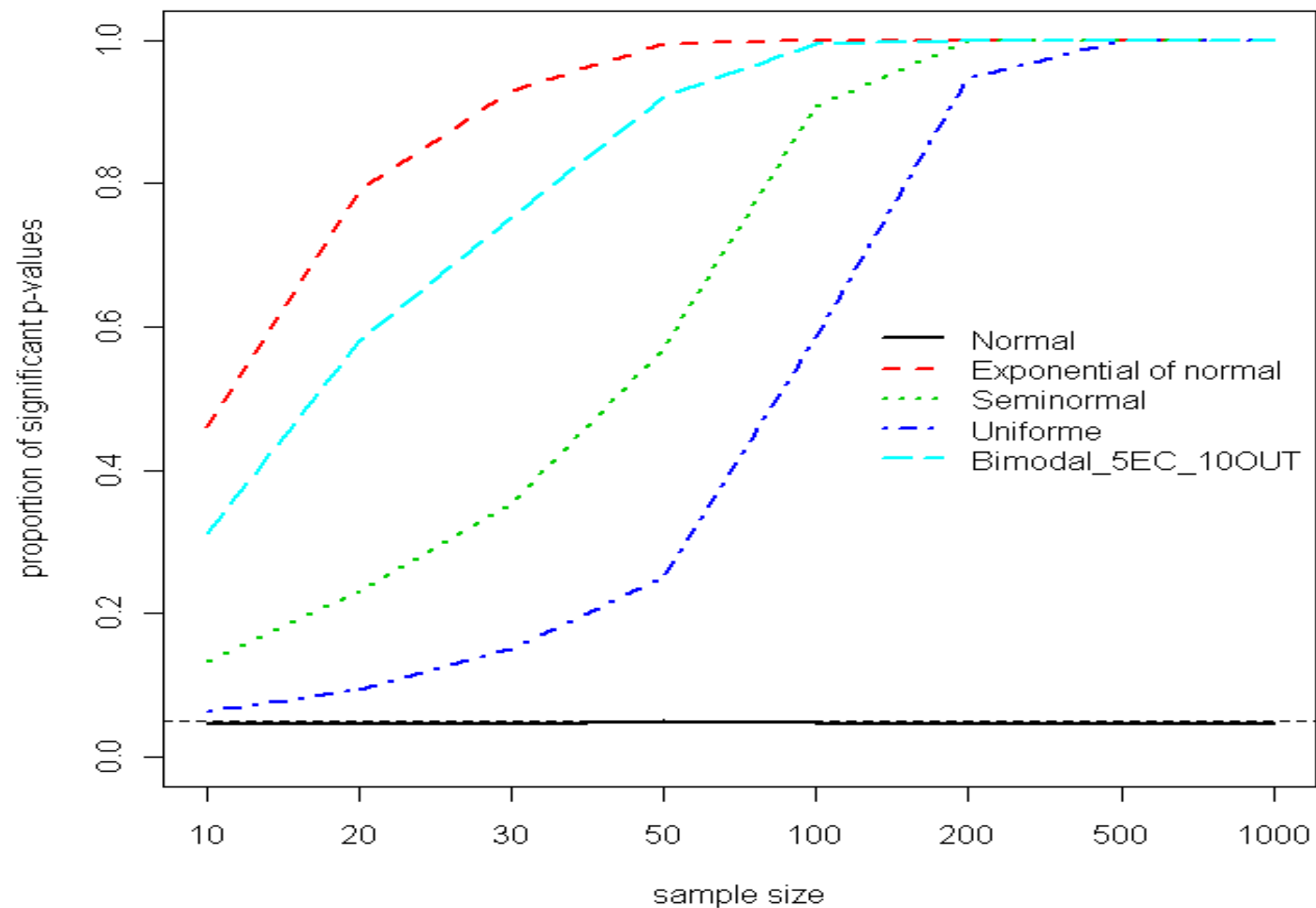
# résumé

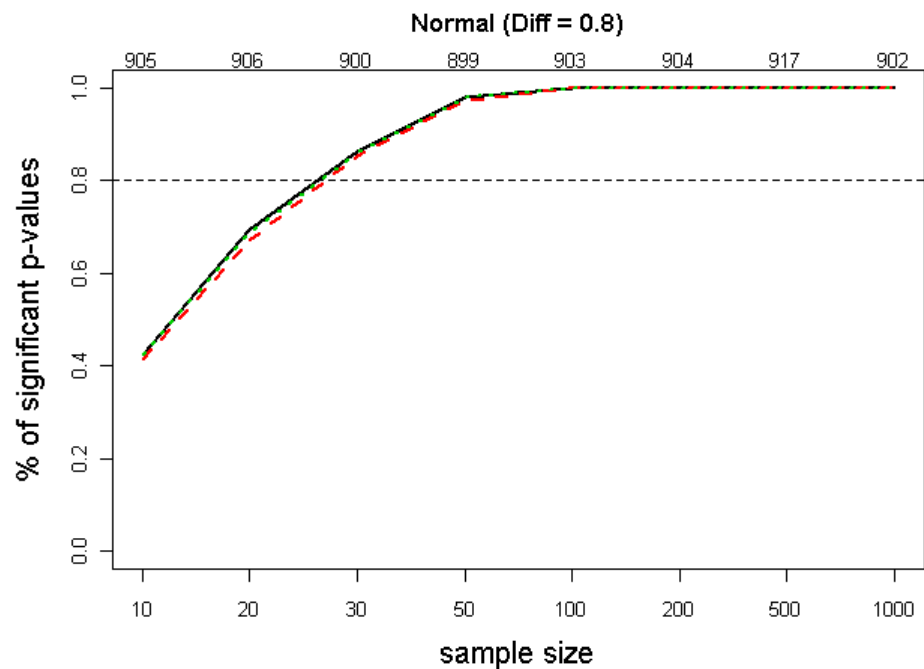
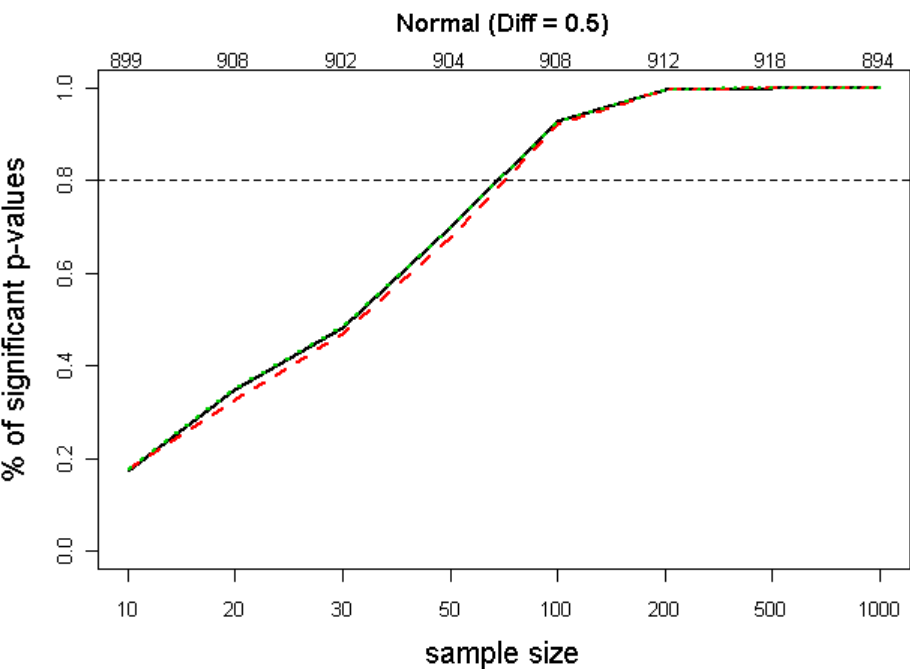
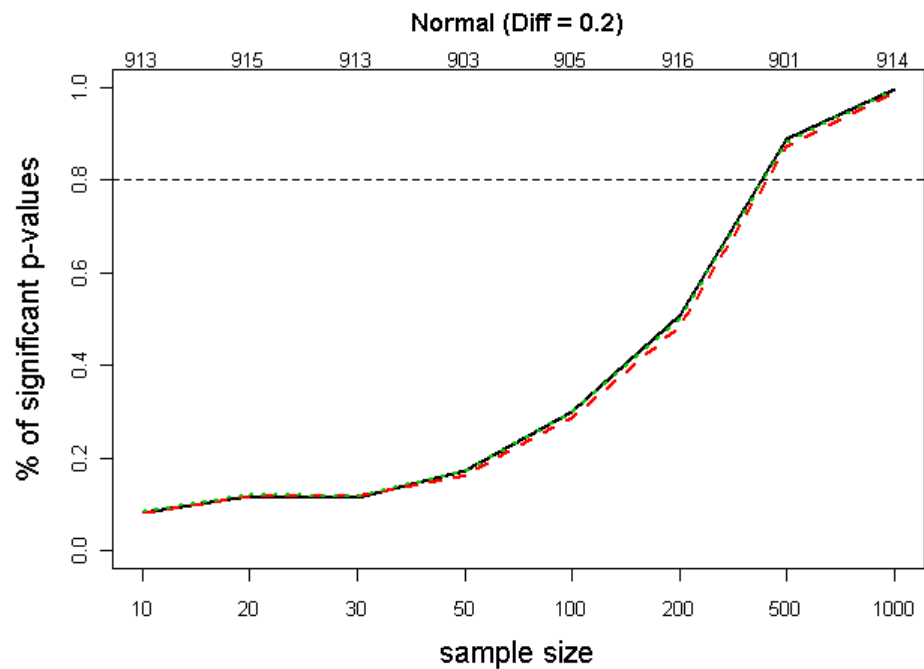
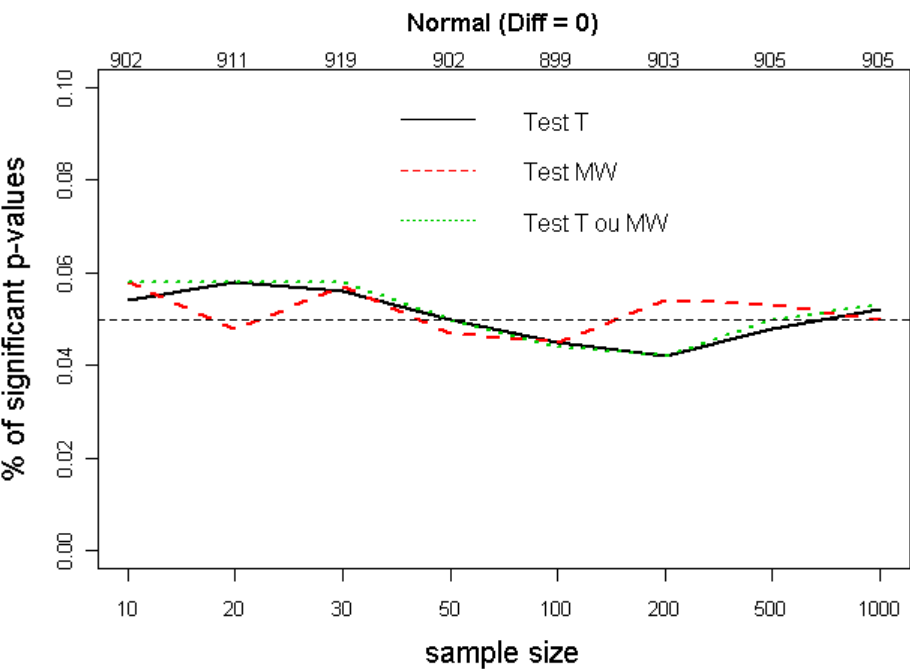
- Sous des conditions de normalité, la concordance entre les deux statistiques est bonne et il y a autant de décisions discordantes « T+ & MW- » que « T- & MW+ »
- Plus on s'éloigne des conditions de normalité en terme d'asymétrie et d'outliers plus on trouve un gain systématique en faveur du MW:
  - Le nombre de cas MW+ et T- augmente
  - Le nombre de cas MW- et T+ diminue

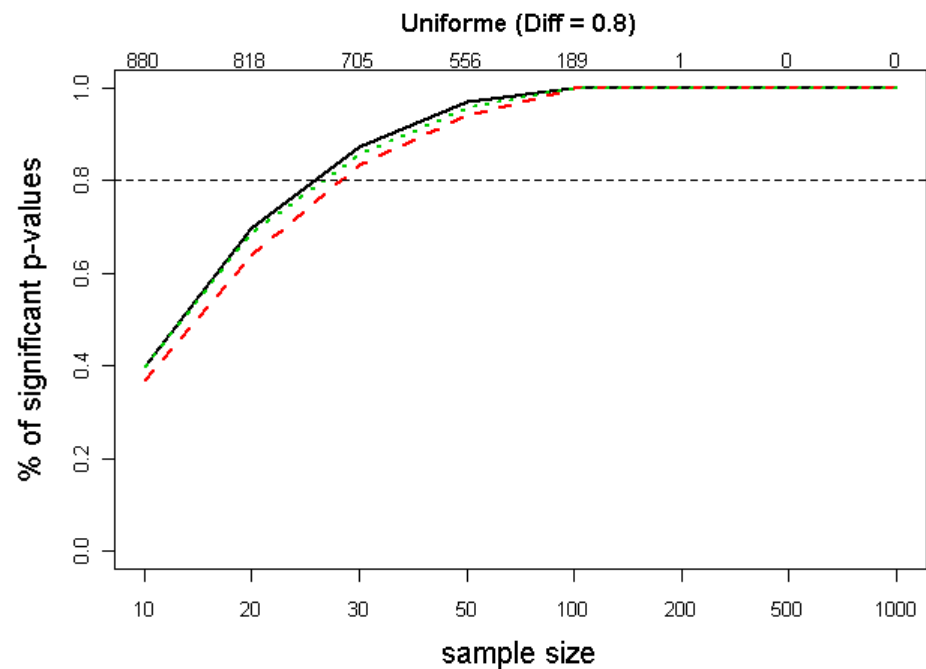
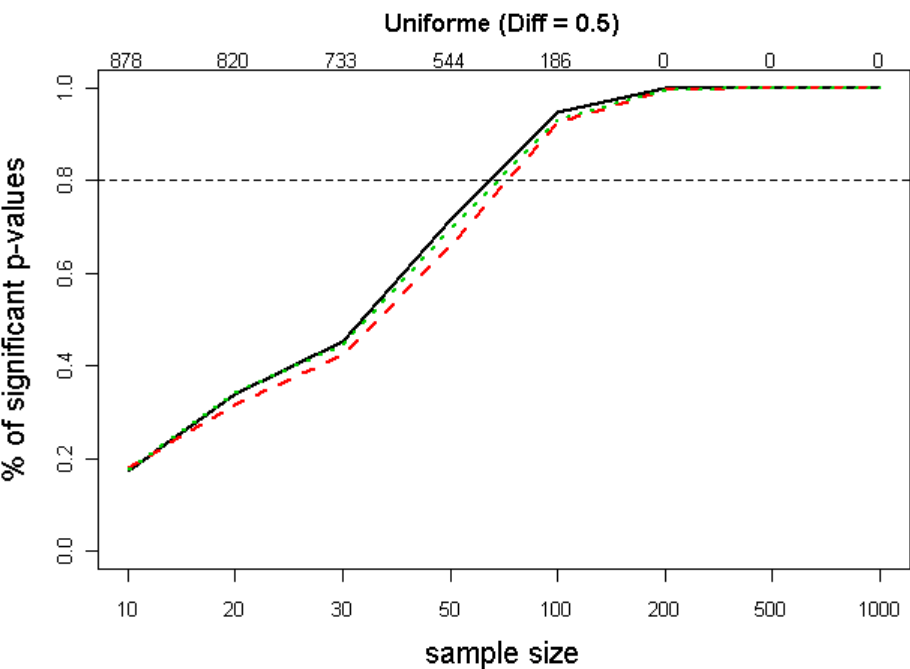
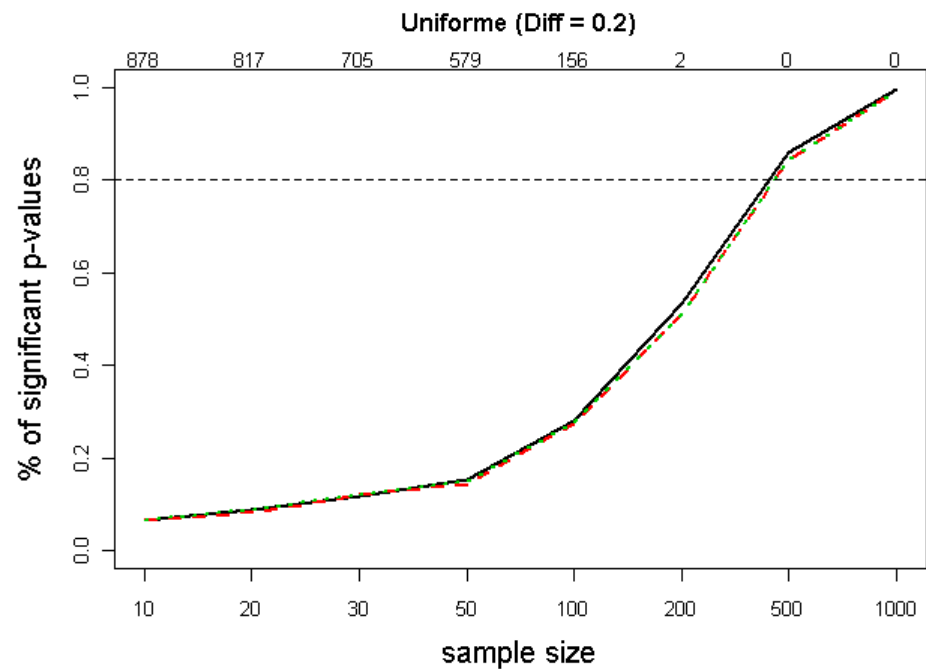
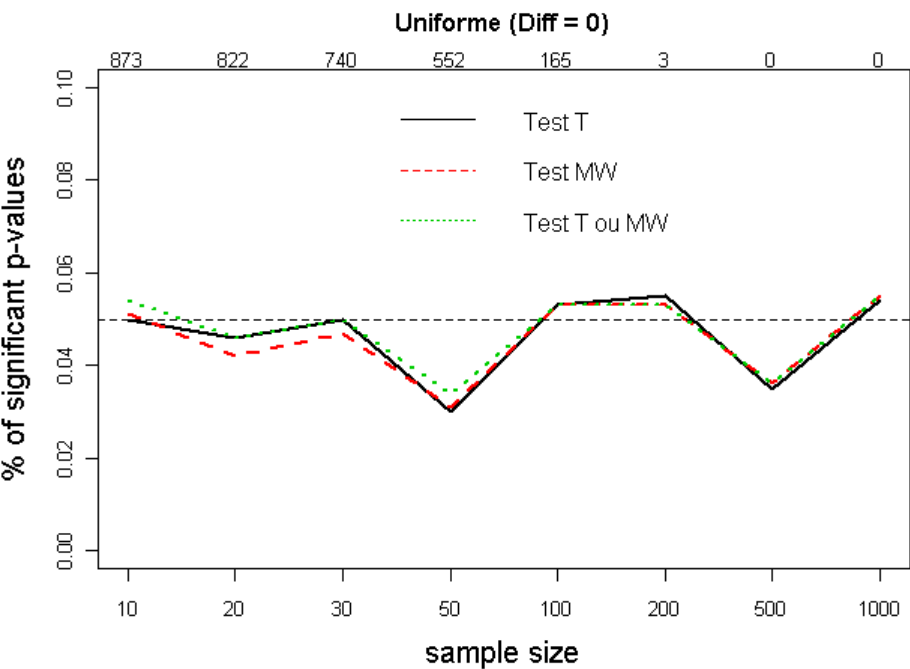
# Analyse de la puissance de 3 stratégies

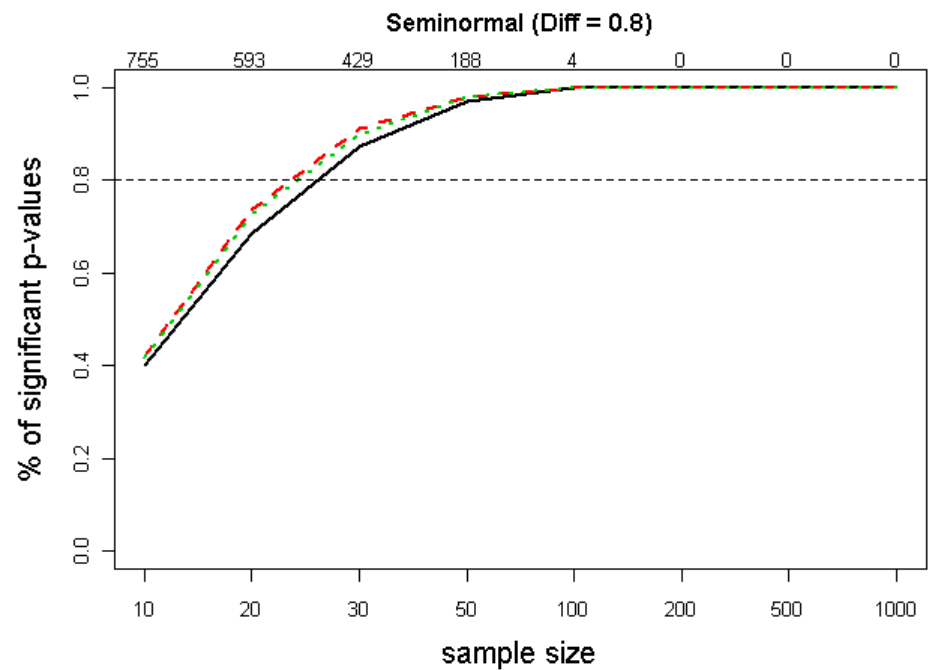
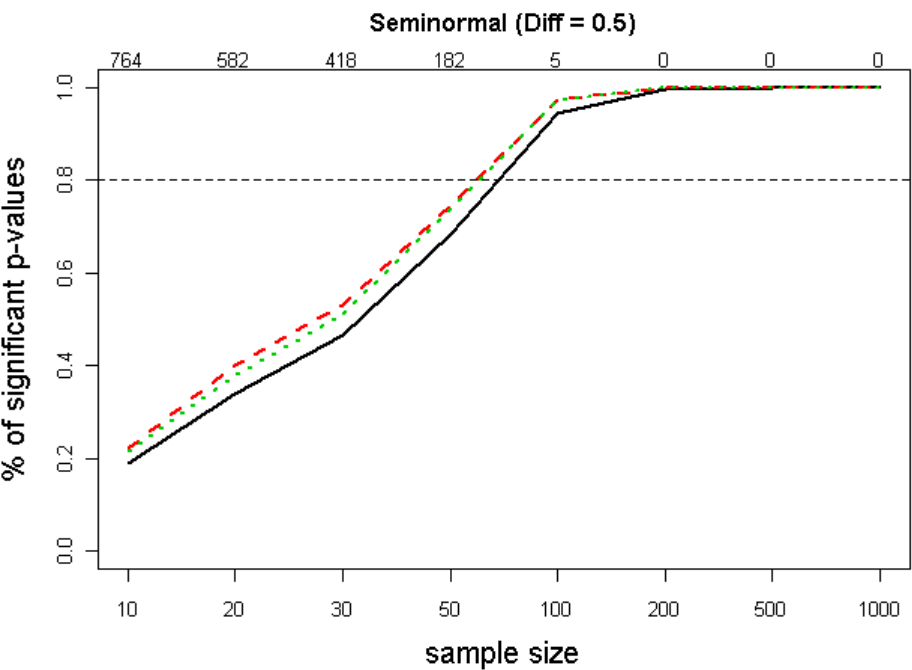
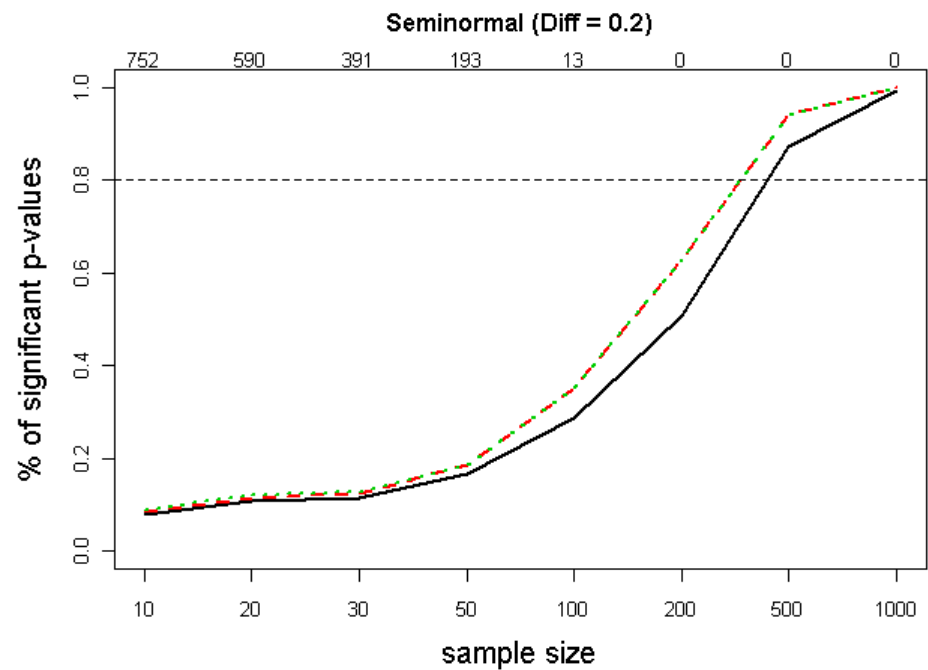
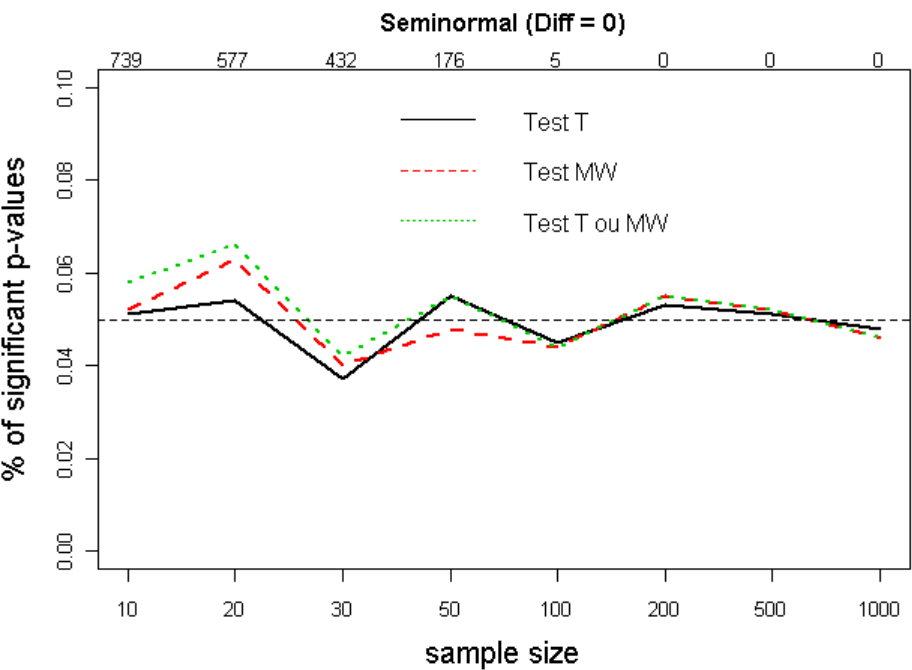
- Stratégie 1: on applique toujours le test t
- Stratégie 2: on applique toujours le test de MW
- Stratégie 3: on applique le test t si le test de normalité (Kolmogorov Smirnov) ne rejette pas  $H_0$ , le test de MW sinon

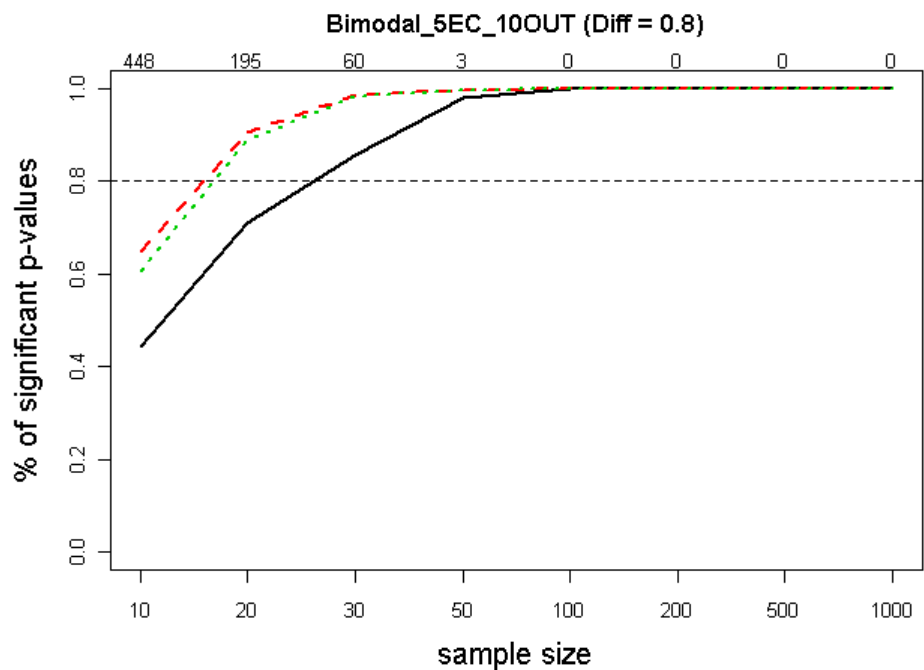
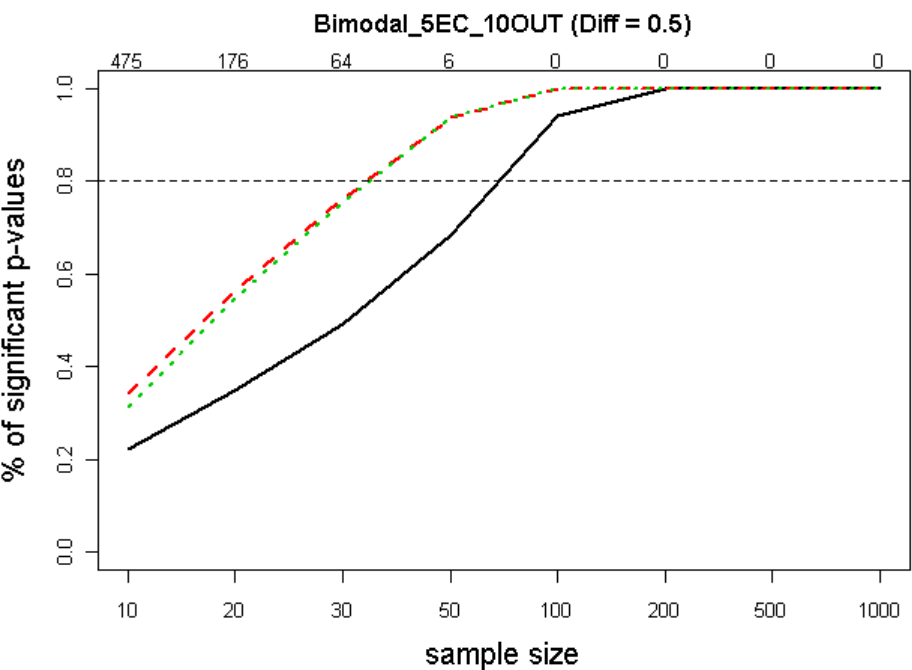
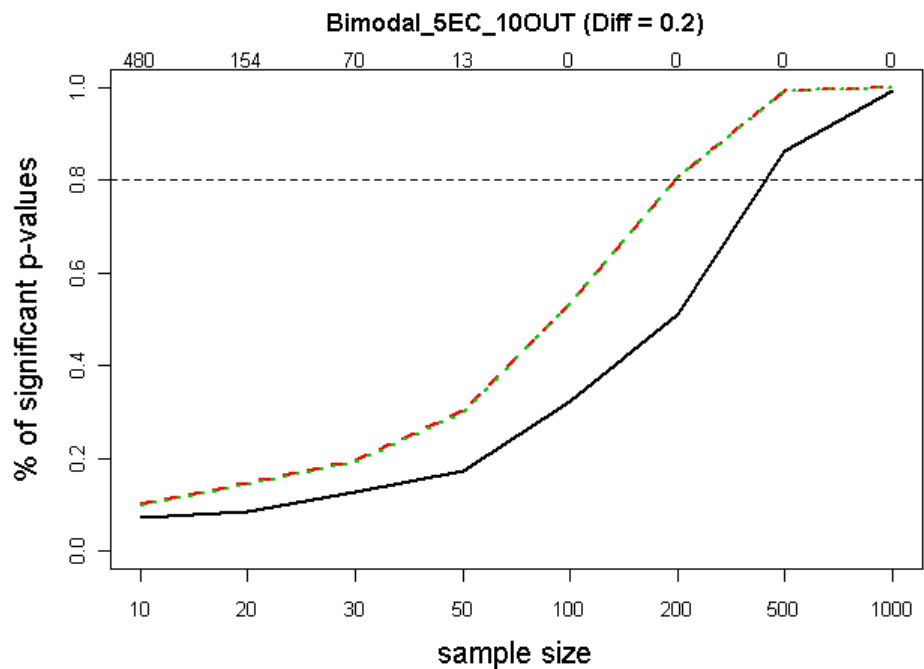
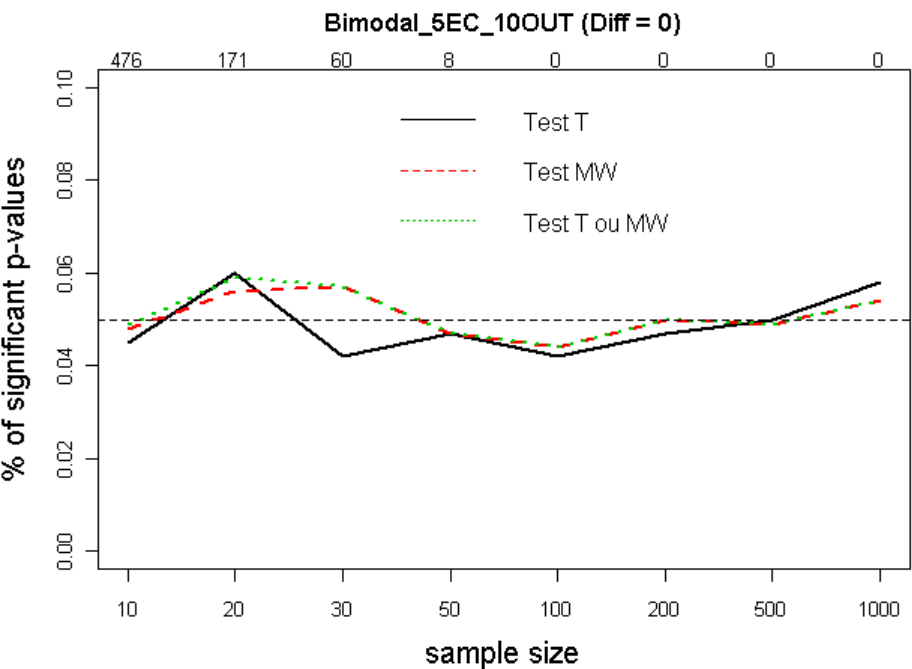
## Kolmogorov-Smirnov test of normality



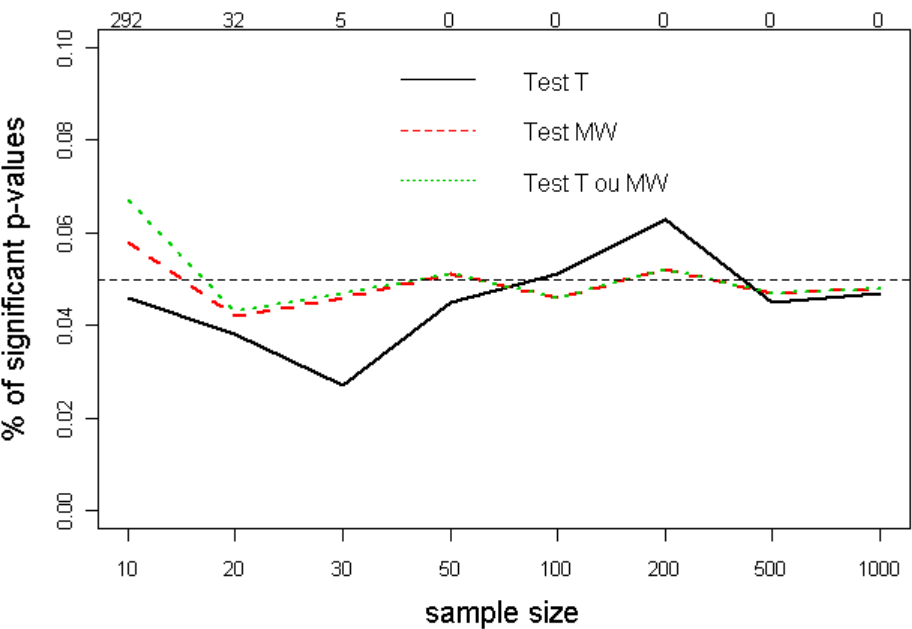




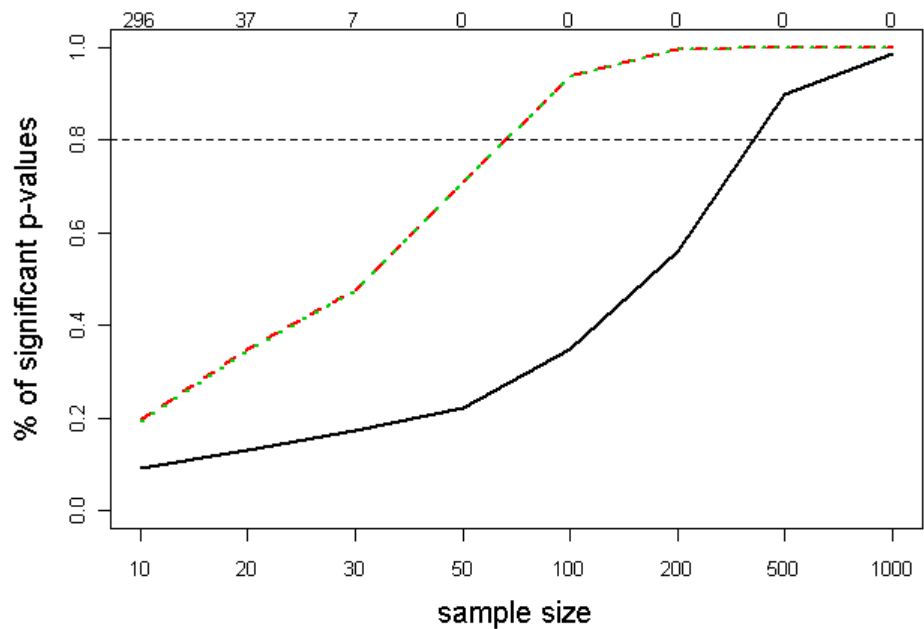




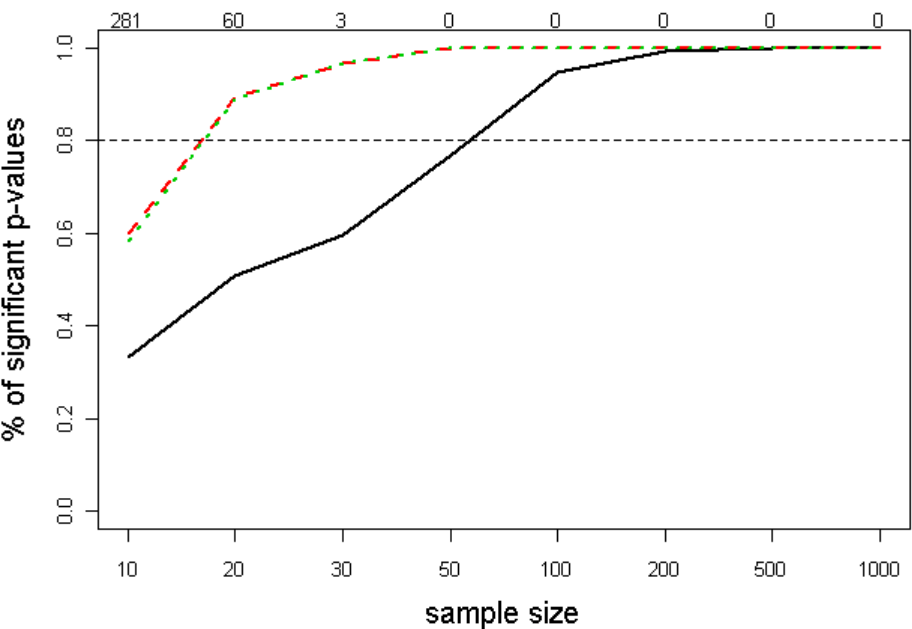
Exponential of normal (Diff = 0)



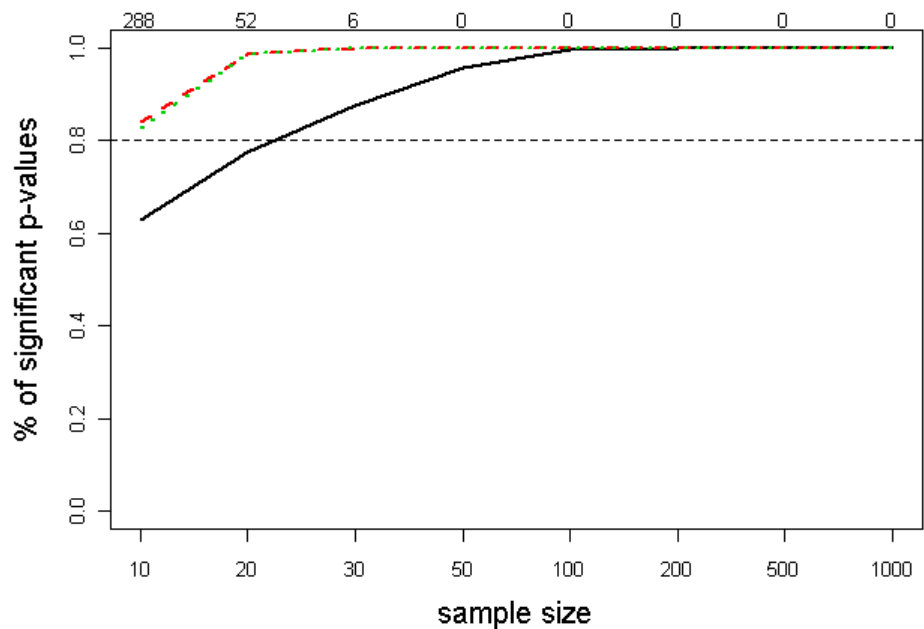
Exponential of normal (Diff = 0.2)



Exponential of normal (Diff = 0.5)



Exponential of normal (Diff = 0.8)





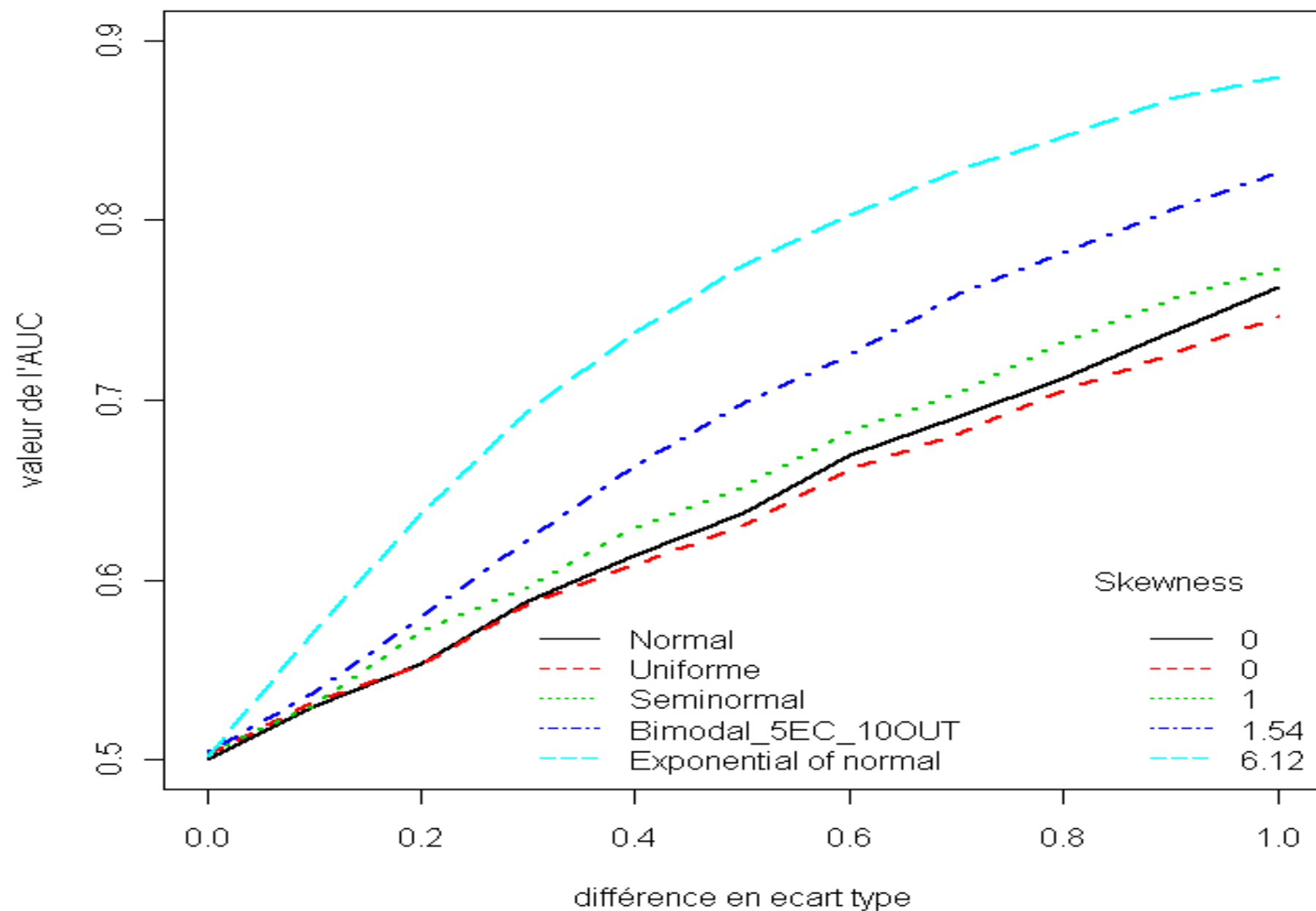
# Conclusion des 3 stratégies

- Le test de normalité n'apporte pas de gain notable !!!
- MW forever...

Conclusion

- Sous distribution symétrique (normale ou uniforme): test t et z ont de meilleures performances que les 3 autres tests mais le gain est très faible
- Sous distribution asymétrique ou en présence d'outliers, gain de puissance du MW et du test robuste pouvant être très important par rapport aux tests t, z et exact
- Pas de « N magique » pour le test de Student (le TCL est trop dépendant de la forme de la distribution)
- Test de normalité inutile

## AUC en fonction de la différence en ET



# Discussion

- Ces résultats confirment certaines publications sur le sujet (sous condition de normalité il faut 0.955 fois moins de sujets avec le test t que le test de MW pour obtenir la même puissance)
- Mais d'autres publications sont en faveur du test t
- Débat ouvert depuis 60 ans !
- On sait que des variances inégales impactent sur les performances du test t et de MW en terme d'erreur de type I. Il s'agit d'une limite de notre étude: on n'a pas fait varier ce paramètre

A guide to choosing an appropriate test

Variiances	Distributions	Sample sizes	<i>t</i> test	WMW test	Welch's <i>U</i> test
Equal	Normal	Equal	*	+	+
		Unequal	*	+	+
	Heavy tailed	Equal	+	*	+
		Unequal	+	*	+
	Skewed	Equal	–	*	–
		Unequal	–	*	–
Unequal	Normal	Equal	+	–	*
		Unequal	–	–	*
	Heavy tailed	Equal	+	–	+
		Unequal	–	–	+
	Skewed <sup>a</sup>	Equal	–	–	–
		Unequal	–	–	–

Symbols: \* = method of choice, best properties, + = acceptable, – = not acceptable.

<sup>a</sup> Transformations are recommended.

Points communs avec notre étude:

- Différentes distributions étudiées

Points différents

- Le ratio des sd varie
- Le ratio des effectifs varie
- La taille de l'effet invariante = 0

Limite de cette étude :

- Le critère de jugement est l'erreur de type I (seulement)

## Misconceptions Leading to Choosing the *t* Test Over the Wilcoxon Mann-Whitney Test for Shift in Location Parameter

Shlomo S. Sawilowsky  
Wayne State University

There exist many misconceptions in choosing the *t* over the Wilcoxon Rank-Sum test when testing for shift. Examples are given in the following three groups: (1) false statement, (2) true premise, but false conclusion, and (3) true statement irrelevant in choosing between the *t* test and the Wilcoxon Rank Sum test.

1. False statement	2. True premise but false conclusion	3. True statement but irrelevant
<p>The Wilcoxon is only for use when the data are originally in the form of ranks</p>	<p>The Wilcoxon's underlying assumptions are weaker (true), therefore the hypothesis being tested is less interesting (false)</p>	<p>The t is a classical test</p>
<p>The Wilcoxon is only for use in the presence of outliers</p>	<p>In terms of central tendency, the Wilcoxon pertains to the median (true), which is less interesting than the mean (false)</p>	<p>Results based on the t have been accumulating for almost a century, permitting direct comparison of results over time</p>
<p>The Wilcoxon should only be used for small samples</p>		<p>The hypotheses being tested for the t and Wilcoxon aren't exactly the same</p>
<p>If a modern procedure should be used, it should be a permutation test, not the W</p>		<p>Even its inventor called the Wilcoxon test a « quick and dirty » or « crude » procedure</p>

# Présentation des tests (4/5)

## test exact

- Test basé sur les permutations
  - Aucune hypothèse sur les distributions si  $n < 50$
  - Approximation normale usuelle (test Z) si  $n \geq 50$
- $H_0$ : les distributions sont identiques dans les 2 groupes

