

Un nouveau test statistique pour la comparaison de proportions

Éric TAILLARD¹, Philippe WAELTI¹, Jacques ZUBER¹

1. EIVD, Route de Cheseaux 1, CP, CH-1401 Yverdon-les-Bains, Suisse
 {eric.taillard, philippe.waelti, jacques.zuber} [at] eivd.ch

Mots-clefs : Test statistique, comparaison d'heuristiques

Qui n'a pas été une fois perplexe en lisant, dans un article comparant deux méthodes d'optimisation, des résultats numériques présentés sous la forme suivante : «Nous avons testé notre méthode A sur un jeu de n exemples de problèmes de la littérature et nous avons réussi à résoudre a de ces exemples. La méthode concurrente B n'a réussi à en résoudre que b , mais elle n'a été testée que sur m exemples de problèmes.»

Un tel résultat peut laisser le lecteur perplexe, car il ne dispose pas immédiatement de la réponse à la question de base qu'il doit légitimement se poser : «Est-ce qu'un taux de succès de a/n est significativement supérieur à un taux de succès de b/m ?» Ce que nous appelons «succès» doit naturellement être défini par l'utilisateur : il peut s'agir de l'obtention d'une solution optimale ou de qualité donnée, ou encore simplement d'une solution admissible. Bien souvent, la réponse à cette question centrale ne peut pas se trouver en appliquant un test statistique standard (basé sur le théorème central-limite) car la taille des échantillons (n et m) est trop petite. En se tournant vers des tests non paramétriques (c.f. [1]), il pourra trouver une variante du test du signe, connue sous le nom de test de Mc Nemar, mais qu'il ne pourra appliquer qu'à des données appariées (i.e. un problème doit être résolu par les deux méthodes).

En optimisation combinatoire, il est cependant fréquent de considérer des jeux de problèmes relativement petits et des tableaux de résultats incomplets (une méthode ou l'autre n'a pas été testée sur certains problèmes). Toutefois, le lecteur sera parfaitement convaincu — à juste titre d'ailleurs — qu'une méthode résolvant les 10 exemples d'un jeu de problèmes est meilleure qu'une méthode ne réussissant à en résoudre aucun. Mais qu'en est-il si le jeu ne comporte que 3 exemples ? On peut montrer qu'un taux de succès de $3/3$ est supérieur à un taux de $0/3$, avec un seuil de confiance supérieur à 98%. Un des buts de cet article est de le montrer.

Soit p_a (respectivement p_b) la probabilité de succès de la méthode A (respectivement, de la méthode B) et n la taille de l'échantillon (nombre d'exécutions) pour la méthode A (respectivement, m pour la méthode B). Pour réaliser un test statistique, on fera l'hypothèse nulle que $p_a = p_b = p$. La valeur de p restant inconnue. Sous cette hypothèse, la probabilité $S(p, a, n, b, m)$ d'observer a succès ou plus pour la méthode A sur n exécutions et b succès ou moins pour la méthode B sur m exécutions est donnée par le produit de deux distributions binomiales :

$$S(p, a, n, b, m) = \left(\sum_{i=a}^n \frac{n!}{i! \cdot (n-i)!} \cdot p^i \cdot (1-p)^{n-i} \right) \cdot \left(\sum_{j=0}^b \frac{m!}{j! \cdot (m-j)!} \cdot p^j \cdot (1-p)^{m-j} \right)$$

Si $\max_{0 < p < 1} S(p, a, n, b, m) \leq 1 - \alpha$, on peut rejeter l'hypothèse nulle avec un niveau de confiance α et accepter l'hypothèse alternative $p_a > p_b$ (i.e. un taux de succès de a/n est plus grand qu'un taux de succès de b/m).

Malheureusement, la valeur de p qui maximise $S(p, a, n, b, m)$ est en général pour le moins difficile à exprimer analytiquement. C'est pourquoi le tableau 1 donne, pour différentes valeurs de n et m et pour un niveau de confiance de 95% les couples (a, b) pour lesquels un taux de succès de a/n est plus élevé qu'un taux de b/m .

Le calcul des niveaux de confiance peut se réaliser *en ligne* à l'URL :

<http://ina.eivd.ch/projects/stamp/>.

Pour de petits échantillons, le nouveau test proposé est plus puissant que le test de Mc Nemar, lorsque ce dernier peut s'appliquer : la probabilité de rejet de l'hypothèse nulle sera plus élevée avec notre test qu'avec le test de Mc Nemar [2].

m	n													
	2	3	4	5	6	7	8	9	10	11	12	13	14	15
2		(3,0)	(4,0)	(5,0)	(5,0)	(6,0)	(7,0)	(7,0)	(8,0)	(9,0)	(9,0)	(10,0)	(11,0)	(11,0)
3	(2,0)	(3,0)	(3,0)	(4,0)	(4,0)	(5,0)	(5,0)	(6,0)	(6,0)	(7,0)	(7,0)	(8,0)	(9,0)	(9,0)
4	(2,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)	(5,0)	(5,0)	(5,0)	(6,0)	(6,0)	(7,0)	(7,0)	(8,0)
5	(2,0)	(2,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)	(4,0)	(5,0)	(5,0)	(5,0)	(6,0)	(6,0)	(7,0)
6	(2,0)	(2,0)	(2,0)	(3,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)	(4,0)	(5,0)	(5,0)	(5,0)	(6,0)
7	(2,0)	(2,0)	(2,0)	(3,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)	(4,0)	(4,0)	(5,0)	(5,0)	(5,0)
8	(2,0)	(2,0)	(2,0)	(2,0)	(3,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)	(4,0)	(4,0)	(5,0)	(5,0)
9	(2,0)	(2,0)	(2,0)	(2,0)	(3,0)	(3,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)	(4,0)	(4,0)	(4,0)
10	(2,0)	(2,0)	(2,0)	(2,0)	(2,0)	(3,0)	(3,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)	(4,0)	(4,0)
11	(2,0)	(2,0)	(2,0)	(2,0)	(2,0)	(2,0)	(3,0)	(3,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)	(4,0)
12	(2,0)	(2,0)	(2,0)	(2,0)	(2,0)	(2,0)	(3,0)	(3,0)	(3,0)	(3,0)	(3,0)	(3,0)	(4,0)	(4,0)

TAB. 1 – Couples (a, b) pour lesquels un taux de succès $\geq a/n$ est significativement plus élevé qu'un taux de b/m , pour un niveau de confiance de 95%.

Remerciements Le présent travail a été en partie financé par la réserve stratégique de la Haute École Spécialisée de Suisse Occidentale (HES-SO), projet LQF01-03.

RÉFÉRENCES

- [1] W. J. Conover, *Practical Nonparametric Statistics*, Wiley, troisième édition, 1999.
- [2] É. D. Taillard, «A Statistical Test for Comparing Success Rates», actes de : *Metaheuristic international conference MIC'03*, Kyoto, Japon, août 2003.