

## Анализ пространства параметров в задачах выбора мультимodelей

А.А. Адуенко<sup>1</sup>, В.В. Стрижов<sup>1,2</sup>

<sup>1</sup>Московский физико-технический институт (государственный университет)

<sup>2</sup>Вычислительный центр им. Дороницына РАН

В работе рассматривается задача выбора мультимodelей при построении modelей в задачах двухклассовой классификации. Мультимodelи являются интерпретируемым обобщением случая одной modelи, позволяющим учитывать неоднородности в данных. Признаковые пространства modelей в мультимodelи могут не совпадать. Кроме того, мультимodelь может содержать большое число близких modelей, что ведет к низкому качеству прогноза и отсутствию интерпретируемости. Для решения этой проблемы предлагается метод статистического сравнения modelей для прореживания мультимodelи. Вводится понятие адекватной мультимodelи, то есть мультимodelи, все modelи в которой являются попарно статистически различимыми.

Для статистического сравнения modelей предлагается ввести функцию близости между апостериорными распределениями параметров modelей. Такая функция должна быть определена для случая пары распределений с несовпадающими носителями, а также не различать два распределения, одно из которых является малоинформативным. Показано, что дивергенция Кульбака-Лейблера, расстояния Дженсона-Шеннона, Хеллингера, Бхаттачарая не удовлетворяют этому требованию. Предлагается функция близости для пары распределений, которая удовлетворяет этим требованиям. Доказаны асимптотические свойства распределения введенной функции близости в условиях истинности гипотезы о совпадении modelей. Использование предлагаемой функции близости распределений для сравнения modelей проиллюстрировано на синтетических данных. С помощью статистических свойств распределения введенной функции близости получены оценки на максимальное количество попарно различимых modelей в мультимodelи для выборки фиксированного размера. Доказана оценка снизу на количество различимых modelей путем построения набора различимых modelей. Для использования в задаче сравнения новой modelи с базовой предложена несимметричная версия введенной функции близости.

## Литература

1. *Bishop C.M.* Pattern recognition and machine learning. // Springer, 2006.
2. *Bishop C.M., Nasrabadi N.M.* Pattern recognition and machine learning. // Journal of electronic imaging, 2007. Vol. 16. No. 4.
3. *Gelman A., Hill J.* Data analysis using regression and multilevel/hierarchical models // Cambridge University Press, 2006.
4. *Siddiqi N.* Credit risk scorecards: developing and implementing intelligent credit scoring // Wiley, 2006.
5. *Hosmer D.W., Lemeshow S.* Applied logistic regression // A Wiley-Interscience Publication, 2000.
6. *Hastie T., Tibshirani R., Friedman J.H.* The Elements of Statistical Learning // Springer, 2001.
7. *Motrenko A., Strijov V., Weber G.W.* Bayesian sample size estimation for logistic regression.
8. *Van den Noortgate W., De Boeck P., Meulders M.* Cross-classification multilevel logistic models in psychometrics // Journal of Educational and Behavioral Statistics, 2003. Vol. 28. No. 4. Pp. 369--386.
9. *Moerbeek M., Van Breukelen G.J.P., Berger M.P.F.* Optimal experimental designs for multilevel logistic models // Journal of the Royal Statistical Society: Series D (The Statistician), 2001. Vol. 50. No. 1. Pp. 17--30.
10. *Link W.A., Barker R.J.* Model weights and the foundations of multimodel inference // Ecology, 2006. Vol. 87. No. 10. Pp. 2626--2635.