

Maximum Mean Discrepancy: A New Analysis Tool for High-Dimensional Biology

Frederick Boehm^{1,*}

May 13, 2021

Abstract

Text of abstract

Contents

1	Specific Aims	2
2	Background and Significance	3
2.1	MMD and its Definition	3
2.2	MMD Properties	3
2.3	Borgwardt et al. [1] Results	3
3	Approach	5
3.1	Approach for Aim 1	5
3.2	Approach for Aim 2	5

¹ University of Massachusetts Medical School

* Correspondence: [Frederick Boehm <frederick.boehm@gmail.com>](mailto:frederick.boehm@gmail.com)

Keywords: keyword 1; keyword 2; keyword 3

Highlights: These are the highlights.

1 Specific Aims

The volume and dimensionality of molecular biology data has exploded over the last twenty years. Following the sequencing of the human genome, technologies such as DNA sequencing, RNA sequencing, and single-cell RNA sequencing have proliferated. The new experimental designs and high dimension of modern biological data require parallel development of new statistical methods and computational tools for their analysis. Maximum mean discrepancy (MMD) is a new statistical tool with novel mathematical properties [5, 4, 3, 1]. We propose to explore its mathematical and statistical properties, to apply MMD methods in the analysis of high-dimensional biological data, and to build user-friendly R software tools that implement MMD methods.

Aim 1: Apply MMD to high-dimensional biological data

Two-sample problems arise in many areas of biomedical research. These include differential gene expression, comparability of data from different laboratories, and classification of cancers. Borgwardt et al. [1] examined performance characteristics of four multivariate two-sample tests in a collection of four biological settings. They found that MMD outperformed the other three multivariate tests.

Aim 2: Create and maintain an open source R package that implements MMD methods

Current R software for MMD methods doesn't provide significance thresholds. Gretton's website hosts Matlab software that implements the methods in Gretton et al. [5], Gretton et al. [3] and Gretton et al. [4]. We'll implement these methods in R and C++ before assembling the code into a well documented, user-friendly R package [7, 2].

2 Background and Significance

MMD is a statistical method for determining whether two high-dimensional samples arise from distinct distributions. Throughout this document, we use the term “sample” to refer to measurements on a collection of subjects from a single group. In introductory statistics courses, one often studies Student’s t-test [8] when discerning whether two samples come from the same distribution. Implicit to the t-test is the assumption that the two samples arise from normal distributions with a single, common variance parameter.

In high-dimensional settings, the analog of Student’s t-test is Hotelling’s T^2 test [6]. Like its univariate counterpart, Hotelling’s test assumes a normal (multivariate joint normal) distribution of the random errors and a shared covariance matrix for the two groups. Departures from these assumptions may give rise to type I and type II errors in testing.

2.1 MMD and its Definition

To combat limitations of previous high-dimensional, two-sample tests, Gretton et al. [4] proposed a kernel-based test for the two-sample problem. Drawing on reproducing kernel Hilbert space (RKHS) theory, they formulate their test as a measure of distance between embeddings of probability measures in a RKHS. Gretton et al. [4] then write the test as a maximization over a subspace \mathcal{F} of their RKHS, \mathcal{H} (Equation (1)).

$$MMD[\mathcal{F}, p, q] = \sup_{f \in \mathcal{F}} (\mathbb{E}_x f(x) - \mathbb{E}_y f(y)) \quad (1)$$

By embedding probability measures in a Hilbert space, and a RKHS in particular, Gretton et al. [5] have the option of using a rich set of analytic tools and theory. They lean heavily on the “kernel trick”, which characterizes the “reproducing” nature of the RKHS \mathcal{H} . The Riesz representation theorem for Hilbert spaces states that there is a feature map, $\phi(x) = k(x, \cdot)$ such that $f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}}$ for any $f \in \mathcal{H}$. Gretton et al. [5] then define an embedding for probability distribution p as an element $\mu_p \in \mathcal{H}$ satisfying $\mathbb{E}_x f = \langle f, \mu_p \rangle_{\mathcal{H}}$.

2.2 MMD Properties

2.3 Borgwardt et al. [1] Results

In aiming to understand and characterize the performance of MMD-based statistical hypothesis tests, we begin by considering the striking results from Borgwardt et al. [1]. Borgwardt et al. [1] obtained two data sets from published breast cancer studies. Both data sets featured gene expression measurements on 2,166 common genes, but the measurements were obtained from two different microarray platforms. Borgwardt et al. [1] emphasized that, if there were no differences between the measurements on the two platforms, then the two data sets could be analyzed together. A joint analysis, when appropriate, would likely lead to greater statistical power to infer meaningful expression signatures for breast cancer. After scaling and centering expression measurements, they examined four tests. For the MMD-based test, they chose a Gaussian kernel with standard deviation set to 20. Borgwardt et al. [1] compared the MMD-based test to standard multivariate tests including Hotelling’s T^2 test, the Friedman-Rafsky multivariate Kolmogorov-Smirnov test, and the Friedman-Rafsky Wald-Wolfowitz test.

In the first of two analyses, the study design consisted of randomly choosing 25 subjects’ vectors of expression measurements from each sample (*i.e.*, each platform). That is, the gene expression vectors of length 2,166 were used. The two resulting subsamples consisted of 25 by 2166 matrices. To these pairs of subsamples, Borgwardt et al. [1] applied the four tests. They repeated the entire procedure 100 times. They recorded the number of instances in which each of the four multivariate tests rejected the null hypothesis that the two subsamples arose from the same platform. Table 1 tabulates the counts for the 100 instances of each test.

The second of two analyses had a similar study design, except that instead of choosing 25 subjects from each of two platforms, Borgwardt et al. [1] first chose one of the two platforms. They then chose two groups of 25 from the chosen platform and compared the resulting subsamples. Again, they repeated the procedure 100 times and recorded the results in Table 1.

Table 1: Microarray cross-platform comparability

platform	null_hypothesis	MMD	t_test	wolfowitz	smirnov
same	accepted	100	100	93	95
same	rejected	0	0	7	5
different	accepted	0	95	0	29
different	rejected	100	5	100	71

The results in Table 1 display perfect accuracy for the MMD-based test: it correctly identified 100 of 100 subsamples chosen from different platforms and 100 of 100 subsamples chosen from the same platform. The two Friedman-Rafsky tests performed at an intermediate level, while the Hotelling’s T^2 test performed poorly. Specifically, the Hotelling’s T^2 test incorrectly called 95 of 100 cases when the two subsamples were chosen from different platforms.

Table 2: Cancer diagnosis

health_status	null_hypothesis	MMD	t_test	wolfowitz	smirnov
same	accepted	100	100	97	98
same	rejected	0	0	3	2
different	accepted	0	100	0	38
different	rejected	100	0	100	62

Table 3: Cancer subtype

subtype	null_hypothesis	MMD	t_test	wolfowitz	smirnov
same	accepted	100	100	95	96
same	rejected	0	0	5	4
different	accepted	0	100	0	22
different	rejected	100	0	100	78

The remarkable performance of the MMD-based tests in these two high-dimensional settings deserves further study. It is from this perspective that we propose research to more completely characterize and understand the MMD-based test.

3 Approach

3.1 Approach for Aim 1

With the goal of enhancing understanding of the MMD-based test in high dimensions, we will undertake the following research plan. First, we'll identify existing R code that implements Hotelling's T^2 test and the two Friedman-Rafsky tests. By achieving **Aim 2**, we'll have R implementations of the MMD-based tests.

3.2 Approach for Aim 2

I'll draw on my extensive experience in statistical software development to create an R package with code in C++ and R. The R package `Rcpp` facilitates the integration of C++ code into R packages. The current code base from Gretton's website is written in Matlab. I'll translate it, line for line, into C++ or R. The proposed R package will have extensive unit tests to ensure fidelity and accuracy of my code and to compare results with those obtained with Gretton's Matlab code. The R package `testthat` makes writing of unit tests especially easy [9].

References

- [1] Karsten M Borgwardt et al. “Integrating structured biological data by kernel maximum mean discrepancy”. In: *Bioinformatics* 22.14 (2006), e49–e57.
- [2] Dirk Eddelbuettel and Romain François. “Repp: Seamless R and C++ Integration”. In: *Journal of Statistical Software* 40.8 (2011), pp. 1–18. DOI: [10.18637/jss.v040.i08](https://doi.org/10.18637/jss.v040.i08). URL: <https://www.jstatsoft.org/v40/i08/>.
- [3] Arthur Gretton et al. “A Fast, Consistent Kernel Two-Sample Test”. In: *NIPS*. Vol. 23. 2009, pp. 673–681.
- [4] Arthur Gretton et al. “A kernel method for the two-sample-problem”. In: *Advances in neural information processing systems* 19 (2006), pp. 513–520.
- [5] Arthur Gretton et al. “A kernel two-sample test”. In: *The Journal of Machine Learning Research* 13.1 (2012), pp. 723–773.
- [6] Harold Hotelling. “The Generalization of Student’s Ratio”. In: *The Annals of Mathematical Statistics* 2.3 (1931), pp. 360–378.
- [7] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria, 2021. URL: <https://www.R-project.org/>.
- [8] Student. “The probable error of a mean”. In: *Biometrika* (1908), pp. 1–25.
- [9] Hadley Wickham. “testthat: Get Started with Testing”. In: *The R Journal* 3 (2011), pp. 5–10. URL: https://journal.r-project.org/archive/2011-1/RJournal_2011-1_Wickham.pdf.