



## **hdmax2, an R package to perform high dimension mediation analysis**

Florence Pittion, Basile Jumentier, Aurélie Nakamura, Johanna Lepeule,  
Olivier Francois, Magali Richard

### **► To cite this version:**

Florence Pittion, Basile Jumentier, Aurélie Nakamura, Johanna Lepeule, Olivier Francois, et al..  
hdmax2, an R package to perform high dimension mediation analysis. 2024. hal-04658960v2

**HAL Id: hal-04658960**

**<https://hal.science/hal-04658960v2>**

Preprint submitted on 3 Sep 2024 (v2), last revised 6 Sep 2024 (v3)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# hdmax2, an R package to perform high dimension mediation analysis

Florence Pittion<sup>1</sup>, Basile Jumentier<sup>2</sup>, Aurélie Nakamura<sup>3</sup>, Johanna Lepeule<sup>3</sup>,  
Olivier François<sup>1,\*</sup> & Magali Richard<sup>1,\*</sup>

<sup>1</sup> Univ. Grenoble Alpes, CNRS, UMR 5525, VetAgro Sup, Grenoble INP, TIMC, 38000 Grenoble, France

<sup>2</sup> Research Center of the Sainte-Justine University Hospital, University of Montreal, Montreal, Quebec, Canada

<sup>3</sup> Université Grenoble-Alpes, Inserm, CNRS, Team of Environmental Epidemiology Applied to Development and Respiratory Health, Institute for Advanced Biosciences, Grenoble, France

\* Co-last authors

**Correspondence:** [olivier.francois@univ-grenoble-alpes.fr](mailto:olivier.francois@univ-grenoble-alpes.fr) & [magali.richard@univ-grenoble-alpes.fr](mailto:magali.richard@univ-grenoble-alpes.fr)

## Abstract

Mediation analysis plays a crucial role in epidemiology, unraveling the intricate pathways through which exposures exert influence on health outcomes. Recent advances in high-throughput sequencing techniques have generated growing interest in applying mediation analysis to explore the causal relationships between patient environmental exposure, molecular features (such as omics data) and various health outcomes. Mediation analysis handling high-dimensional mediators raise a number of statistical challenges. Despite the emergence of numerous methods designed to tackle these challenges, the majority are limited to continuous outcomes. Furthermore, these advanced statistical approaches have yet to find widespread adoption among epidemiologists and health data scientists in their day-to-day practices. To address this gap, we introduce a method specifically tailored for high-dimensional mediation analysis using the max-squared method (HDMAX2). This tool aims to bridge the current divide by providing a practical solution for researchers and practitioners eager to explore intricate causal relationships in health data involving complex molecular features. Here we improve the HDMAX2 method, and expand its capabilities to accommodate multivariate exposure and non-continuous outcomes. This improvement enables its application to a diverse array of mediation analysis scenarios, mirroring the complexity often encountered in healthcare data. To enhance accessibility for users with varying expertise, we release an R package called `hdmax2`. This package allows users to estimate the indirect effects of mediators, calculate the overall indirect effect of mediators, and facilitates the execution of high-dimensional mediation analysis.

**Keywords:** High-dimension; Mediation; Multivariate analysis; Confounding effect; Causal analysis

## Introduction

When a statistical association is observed between an external exposure ( $X$ ) and an individual outcome ( $Y$ ), one or more intermediate variables ( $M$ ) (such as gene expression or epigenetic changes) may mediate this effect. Identifying and assessing the mediating role of these variables in the effect of  $X$  on  $Y$  is crucial for deciphering underlying causal mechanisms in epidemiological and clinical research. This process, known as mediation analysis, involves studying mediator variables to define the causal structure between  $X$  and  $Y$ . The mediated effect, termed the indirect effect, is equal to the portion of the effect of  $X$  on  $Y$  mediated through  $M$ , to distinguish from the direct effect of  $X$  on  $Y$  unexplained by  $M$  (Richiardi et al., 2013). Historically, mediation analysis has predominantly focused on univariate mediation (Baron and Kenny, 1986), running separate statistical tests for the effects of  $X$  on  $M$  and  $M$  on  $Y$ , followed by estimation of the indirect effect (Imai et al., 2010; Sobel, 1982). However, in the realm of high-dimensional molecular data ( $M$ , e.g., omic data), extending mediation analysis to high dimensions poses challenges, including correction for multiple testing, controlling the false discovery rate (FDR), addressing reverse causation, adjusting for confounding effects, considering interactions among mediators, and integrating multimodal data types (Blum et al., 2020; Zeng et al., 2021). Currently, there remains no consensus on the optimal combination of models and methodologies for high-dimensional mediation analysis. With the increasingly prevalent use of next-generation sequencing technologies, there is now an urgent need to develop high-performing methods in high dimensionality and make them accessible. While recent methods have partially addressed these needs (Dai et al., 2022; Djordjilović et al., 2022; Sampson et al., 2018; H Zhang, Zheng, Z Zhang, et al., 2016), particularly in controlling the type I error in high dimension, they do not account for unmeasured confounding factors and do not allow for the consideration of multiple exposures.

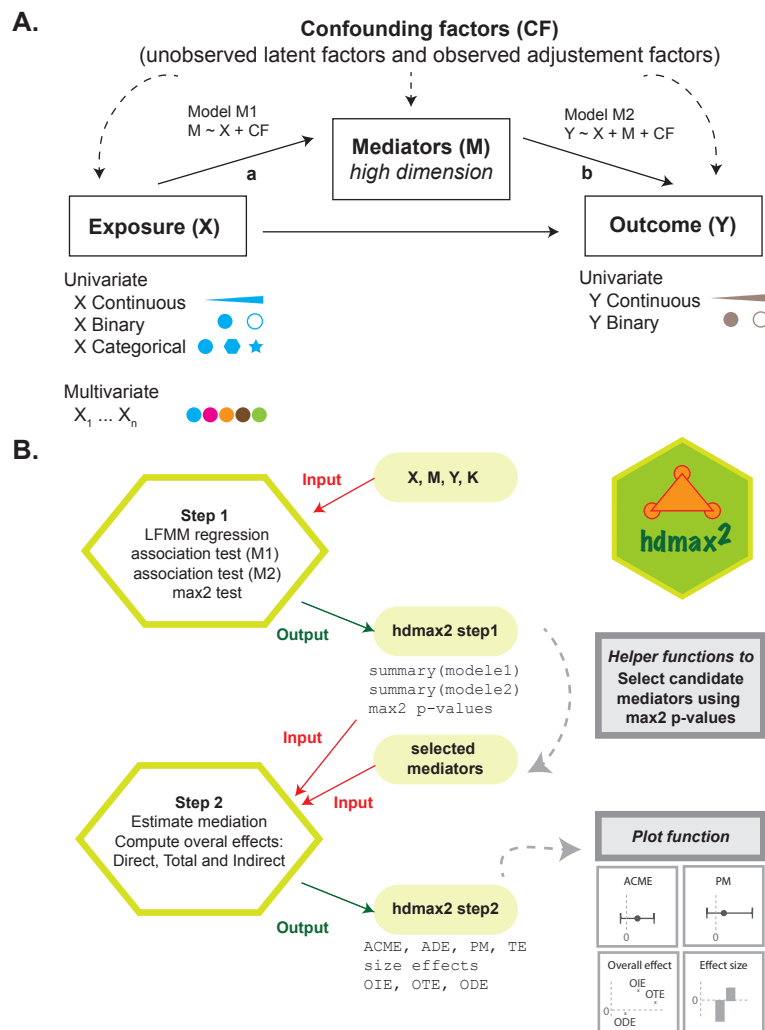
In this article, we introduce an R package called `hdmax2`. The HDMAX2 method was originally proposed by Jumentier et al. (Jumentier et al., 2023). The fundamental concept behind HDMAX2 methods is to use a latent factor mixed regression model for estimating unobserved latent factors while conducting high-dimensional association analysis. HDMAX2 also implements a novel procedure known as the max-squared test to assess the statistical significance of potential mediators. Finally, HDMAX2 enables the calculation of an overall indirect effect from a single model that includes all mediators simultaneously. This approach therefore takes into account correlations between mediators. A case study investigating the impact of maternal smoking on birth weight identified epigenetic regions mediating the indirect effect of this exposure (Jumentier et al., 2023). In this article, we introduce an enhanced version of the HDMAX2 method, expanding its functionalities and making the method accessible to practitioners through a packaged format. The `hdmax2` program has been enriched with numerous features, including the ability to accommodate various types of variables in the exposure (continuous, binary, categorical, and multivariate), as well as the capability to incorporate binary outcomes. This package enables users to (i) investigate associations between the variables  $X$ ,  $M$ , and  $Y$ , (ii) compute the mediated effect for each potential mediator, (iii) assess the overall indirect effect for the total model, and (iv) visualize these results. A graphical model of the package is presented in Fig 1. The package is open-source and accessible on our GitHub page at <https://github.com/bcm-uga/hdmax2>.

## Materials and methods

The package comprises a set of core functions along with a visualization function. Its usage is guided through a vignette and test datasets [https://bcm-uga.github.io/hdmax2/articles/hdmax2\\_tutorial.html](https://bcm-uga.github.io/hdmax2/articles/hdmax2_tutorial.html). The statistical methodologies embedded within the package are described below, with the initial application case documented by Jumentier et al., 2023.

**Input data.** The `hdmax2` package is designed to accept one or several exposure variables,  $X$ , which can be continuous, binary, or categorical. The user must provide an exposure *data.frame* having at least one column

**Figure 1. Graphical model of the `hdmax2` R package.** **A.** Acyclic graph of HDMAX2 method for conducting high-dimension mediation analysis. Exposure variables ( $X$ ) can be univariate or multivariate, the colors represent different types of variable, and the symbols schematically represent the modalities of the same variable. Intermediary variables ( $M$ ) are continuous variables. The outcome variable ( $Y$ ) can be binary or continuous. Confounding factors represent observed adjustment factors and unobserved latent factors estimated by LFMM (Latent Factors mixed model) regression. The values  $a$  and  $b$  represent the effect sizes for each regression of the mediation model. **B.** Core structure of the `hdmax2` package. Step 1: First, latent factors are estimated by LFMM multivariate regression. Then, the association of  $X$  and  $M$  and of  $M$  and  $Y$  are statistically tested in simple regression models, and the significance values obtained are combined, for each potential mediator. Step 2: the indirect effects of the mediator variables are estimated, along with various causal measures of interest.



as input. In the R language, categorical variables are encoded as factor objects. The function `as.factor()` can be used to encode categorical variables. The functions `levels()` and `ordered()` can be used to define the order of the modalities of categorical variables. By convention, `hdmax2` uses the first modality as a reference to calculate the effects associated with the other modalities of the variable, as encoded by default in `lm()` function in R. Continuous intermediary variables, denoted as  $M$ , are represented as a matrix encompassing potential mediators, such as methylome or transcriptome molecular features. The matrix  $M$  should be entered as a separated input, without missing values. The outcome variable, denoted as  $Y$ , corresponds to a vector, which supports both continuous and binary formats. Continuous and binary variables must be

encoded in numeric format. Optional covariates,  $Z$ , can be included as observed adjustment factors in the model. The package `hdmax2` also takes as input the number of latent factors to be estimated,  $K$ . The number of latent factors is usually estimated by applying PCA on the potential mediators matrix and using elbow curve approach on screeplot.

**Step 1: Identification of potential mediators.** The function `hdmax2::run_AS()` evaluates the association between exposure variables, intermediary variables and the outcome variable using a latent factor mixed model (LFMM Caye et al., 2019) to estimate  $K$  unobserved latent factors  $U$ . First this function tests the significance of association between the exposure variables and the potential mediator variables. Then it tests association between the potential mediator variables and the outcome variable. Finally it evaluates the significance of the indirect effects.

For univariate (continuous, binary and categorical) exposure, a significance value  $P_{1,j}$  is computed by the `hdmax2` program for the test of a null effect size for exposure variable  $X$  on intermediary variable  $M_j$ , for each  $j$ . For multivariate exposure, the `hdmax2` program applies partial regression models, and returns a single  $P$ -value,  $P_{1,j}$ , for the test of a null effect size of the full exposure variable  $X$  on the intermediary variable  $M_j$ , for each  $j$  (see Fig 1.A, regression 1). Then, the `hdmax2` program returns a significance value  $P_{2,j}$  corresponding to the association of each intermediary variable  $M_j$  with the outcome variable  $Y$  (see Fig 1.A, regression 2).

Finally `hdmax2` identifies potential mediators by combining the significance values  $P_1$  and  $P_2$  to compute a  $P$ -value for each intermediary variable using the max-squared ( $\max^2$ ) test (equation (1)). This test rejects the null hypothesis of no effect of exposure on potential mediators or no effect of potential mediators on the outcome.

$$P = \max(P_1, P_2)^2. \quad (1)$$

The `hdmax2::run_AS()` function returns an object of class `hdmax2_step1`, including the following attributes:

- the mediation  $P$ -values resulting of the max-squared test,
- $U$ , score matrix for the  $K$  unobserved latent factors calculated from an LFMM regression (model M1),
- the input variables of the model: exposure variables  $X$ , the outcome variable  $Y$ , and adjustment factors  $Z$  (when applicable).

**Step 2: Estimation of indirect effects.** The function `hdmax2::estimate_effect()` takes as input an object `hdmax2_step1` and a list of potential mediators  $M^S$  to be analyzed in subsequent steps. The subset  $M^S$  is defined by the user based on the output of `hdmax2_step1`. We provide a series of helper functions to guide the user in the selection of  $M^S$ . These functions include a False Discovery Rate (FDR) control approach, or the possibility to aggregate mediator regions according to their location on the genome. Illustrations of these approaches can be found in Jumentier et al., 2023 and in the `hdmax2` vignettes.

For each univariate exposure variable and the subset of mediators  $M^S$ , the `hdmax2::estimate_effect()` function computes several estimates to evaluate the indirect effects in the path between exposure variables and the outcome variable. Initially, this function assesses each mediator variable  $M_j^S$  individually and computes causal measures of interest such as (i) the Average Causal Mediated Effect (ACME, corresponding to the indirect effect) and (ii) the Proportion Mediated (PM). The ACME differs from the Average Direct Effect (ADE), which represents the unmediated effect. PM corresponds to the proportion of the total effect that is mediated by the mediator (ratio of the indirect effect to the total effect). ACME and PM are computed by the `mediation::mediate()` function of the package `mediation`, that automatically detects the type of statistical model used in the mediation analysis (Tingley et al., 2014). The function `mediation::mediate()` calculates uncertainty estimates by a quasi-Bayesian Monte Carlo approach described in Imai et al., 2010. In addition, it estimates the intermediary effect sizes  $a_j$  and  $b_j$  and their standard deviations (see Fig 1A). Eventually, `hdmax2`

calculates an Overall Indirect Effect (OIE) from a single model that includes all mediators  $M^S$  simultaneously. The OIE corresponds to the sum of the indirect effect associated with all mediators. The confidence interval (CI) of the OIE is estimated by a bootstrap approach. Along with the OIE, `hdmax2` estimates the Overall Total Effect (OTE) corresponding to the effect of exposure variables on the outcome variable, as well as the Overall Direct Effect (ODE) corresponding to the effect of exposure variables on the outcome variable when the mediators  $M^S$  are included in the model. For categorical variables, all estimates (ACME, PE, OIE, OTE, ODE and size effects) are calculated relative to a reference corresponding to the first modality of the variable. In the case of a multivariate exposure, each variable is treated independently, the other variables being included in the covariable matrix of the mediation model.

The `hdmax2::estimate_effect()` function returns an object of class `hdmax2_step2`, including the following attributes:

- ACME (average causal mediated effect), ADE (average direct effect), PM (proportion mediated) and TE (total effect), for each mediator,
- OIE (overall indirect effect), OTE (overall total effect) and ODE (overall direct effect),
- summaries of regression models adjusted during the mediation analysis.

The function `hdmax2::plot()` takes as input an object `hdmax2_step2` and enables graphical visualization of mediated effects. This function returns an ACME forest plot, a PM forest plot, a plot of the overall effects and a plot of the indirect effect sizes  $a$  and  $b$ .

**Data collection and preprocessing.** The breast cancer dataset was collected from the TCGA-BRCA public repository. DNA methylation data underwent filtering to remove probes containing NA values, resulting in the retention of approximately 20,000 CpG sites. Gene expressions were normalized using standard DESeq2 parameters (Love et al., 2014) and pseudo-log transformed. The multiple sclerosis dataset was collected from the GEO public repository (accession number : GSE137143 (Kim et al., 2021)). RNA-seq data underwent normalization using standard DESeq2 parameters (Love et al., 2014) and were filtered to retain only coding genes with detectable expression ( $>0$  in at least one sample). Filtered data were subsequently transformed using pseudo-logarithmic transformation.

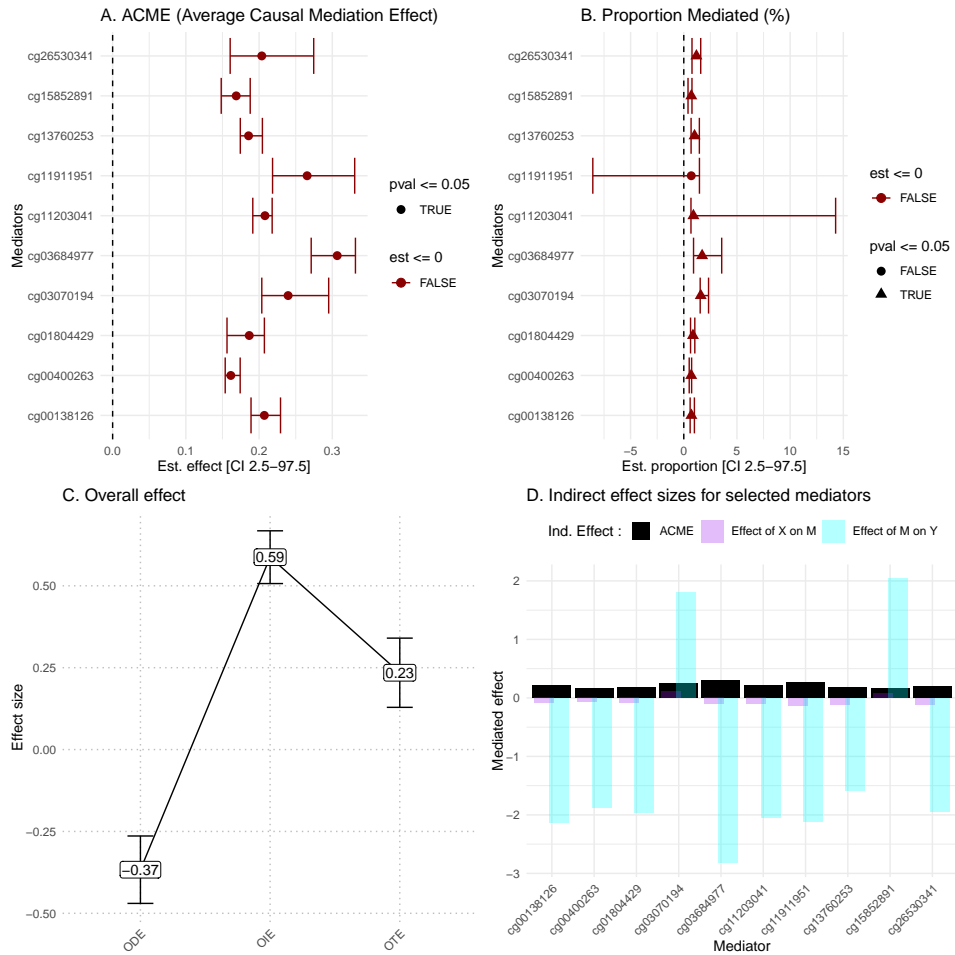
**Bioinformatic analysis.** Gene set enrichment analysis were performed using the `fgsea` (Korotkevich et al., 2021) and `msigdbR` (Dolgalev, 2024) R packages, using the defaults parameters. Gene ranks correspond to  $-\log_{10}(\text{max-squared pvalues})$ .

## Results

**Description of the package vignette** We propose several vignettes to explain the use of the package based on simulated data: an analysis case with a univariate exposure, an analysis case with a multivariate exposure and a vignette to illustrate the usage of helper functions.

**First use case : HER2 and breast cancer** In this example study, we employed mediation analysis to assess the potential causal role of DNA methylation in the pathway linking HER2 status of Breast Cancer to a survival prognostic factor, namely the risk score. Our investigation utilized data from the TCGA-BRCA repository. The risk score was derived from a six-gene expression signature, as described by (Yin et al., 2024), and is inversely correlated with patient survival. HER2 positive status contributes to the molecular subtyping of breast cancer, which includes 'estrogen and progesterone receptor-positive, HER2 negative' (Luminal A), 'hormone positive and HER2 negative' (Luminal B), as well as 'HER2 positive' and 'triple-negative breast cancer (TNBC)' subtypes (Orrantia-Borunda et al., 2022). The computation of the risk score was based on RNAseq data, following the methodology outlined by (Yin et al., 2024). We categorized patients based on HER2 status, dividing them into HER2-positive ( $n=176$ ) and HER2-negative ( $n=752$ ) groups, while filtering out equivocal ( $n=22$ ) and indeterminate ( $n=8$ ) cases. Our analytical model included  $X$  as a binary exposure variable representing HER2 expression

**Figure 2. Summary of the breast cancer use case. A.** Estimates of indirect effect (ACME) and **B.** proportions of mediated effect (PM) for the top 10 mediators. The effect estimate is represented by a dot and its 95% CI by the bar. Symbols correspond to the significance cut off of 5% (square for p-value  $\geq 0.05$ , circle p-value  $< 0.05$ ). Colors correspond to the sign of the effect (green for estimated effect  $\leq 0$ , red for estimated effect  $> 0$ ). **C.** Effect sizes of Overall Direct Effect (ODE), Overall Indirect Effect (OIE) and Overall Total Effect (OTE). Error bars correspond to standard deviation (ODE and OTE) or confidence interval (OIE). **D.** Indirect effect sizes for the selected mediators. Black corresponds to the ACME, violet to the effect of exposure  $X$  on mediator  $M$  in the model  $X \sim M$ , and blue corresponds to the effect of mediator  $M$  on outcome  $Y$  in the model  $Y \sim M + X$ .



(0 = negative, 1 = positive),  $M$  encompassing 20,000 methylation probes, and  $Y$  denoting the continuous risk score.

In our analysis, after adjusting for the confounding effect of age, we found that the total effect of HER2-positive status resulted in a 0.30 higher risk score (t-test,  $p=0.007$ ,  $sd = 0.11$ ). For the initial step of the HDMAX2 approach, we opted to use  $K = 2$  latent factors in the association study. Subsequently, we identified the top 10 potential mediators with the lowest max-squared p-values. We then estimated the individual indirect effect of each mediator by computing the Average Causal Mediation Effect (ACME) as depicted in Fig 2A, along with the proportion mediated shown in Fig 2B. Following this, we calculated the Overall Indirect Effect and its corresponding Overall Direct Effect, as illustrated in Fig 2C. In our analysis, we observed a negative direct effect of HER2 on the risk score, suggesting that HER2 expression in breast tumors has a protective effect on survival. However, the indirect effect mediated by the top 10 CpG probes corresponded to a 0.59 increase in the risk score (standard deviation = 0.07). This indicates that the mediated effect is detrimental to patient survival,



resulting in the observed total effect of 0.30 on the risk score. Interestingly, in our analysis, we noted that for 8 out of the 10 identified mediators, the positive Average Causal Mediation Effect (ACME) resulted from a double negative effect: firstly, from HER2 status ( $X$ ) to methylation ( $M$ ), and secondly, from methylation ( $M$ ) to risk score ( $Y$ ) (Fig 2D). This observation suggests a complex interplay between the exposure, mediators, and outcome, where the presence of the mediator counterbalance the negative effect of the exposure on the outcome. This phenomenon highlights the complex relationships within biological pathways involved in tumor growth and patient survival. This result highlights the importance of considering mediators to understand the mechanisms underlying seemingly simple observed associations. Remarkably, most of the top 10 identified mediators were associated with genes known to be involved in breast cancer biology, thus supporting the biological relevance of our approach (see Table 1).

**Table 1.** Top10 mediators. ID corresponds to the CpG probe name. Chromosome and Start are coordinates provided by the Illumina services. Gene Symbol corresponds to genes known to be associated with the CpG probes (Illumina annotation file). Pubmed hits corresponds to the number of output from the search "(Breast cancer) AND ('Gene Symbol')".

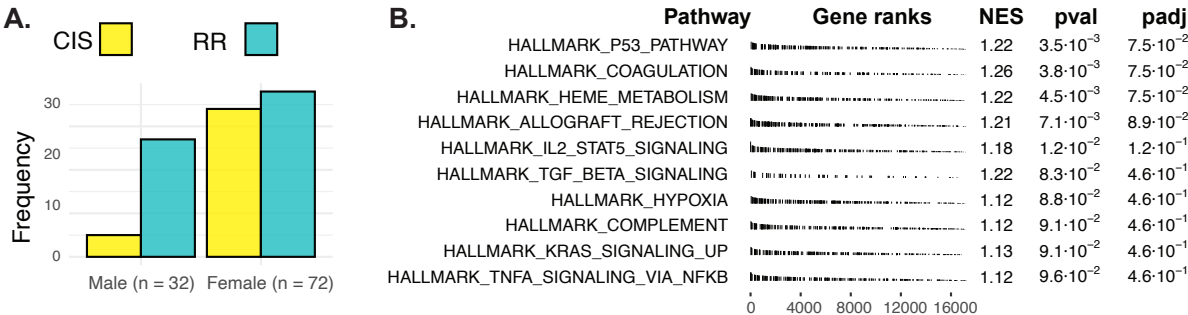
ID	Chromosome	Start	Gene Symbol	Pubmed hits
cg11911951	chr16	1380215	AL031721.1;UNKL	0;0
cg26530341	chr8	23225840	RP11-1149O23.3;TNFRSF10A	0;54
cg01804429	chr4	68350359	YTHDC1	7
cg03684977	chr17	39737550	GRB7	123
cg11203041	chr12	16347252	MGST1	6
cg00138126	chr20	57711641	PMEPA1;RP5-1059L7.1	17;0
cg13760253	chr8	66020465	DNAJC5B	0
cg03070194	chr1	109668062	GSTM2;GSTM4	9;7
cg15852891	chr5	77639121	OTP	8
cg00400263	chr20	59939146	FAM217B;PPP1R3D	0;0

**Second use case : Gender and multiple sclerosis subtypes** In this second case study, we conducted mediation analysis to explore the role of gene expression in the pathway linking patient gender to the occurrence of Multiple Sclerosis (MS) specific subtypes : Clinically Isolated Syndrome (CIS) and Relapsing-Remitting (RR). We used the publicly available dataset GSE137143 (Kim et al., 2021). This dataset comprises peripheral blood mononuclear cells (PBMCs) from healthy individuals and patients with MS. Upon observing a significant decrease in CIS-MS occurrence among women (see Fig 3A), we sought to investigate this phenomenon further. Although the prevalence of MS is known to be higher in women (Harbo et al., 2013), men tends indeed to present worst disease progression (Voskuhl et al., 2020). We found that being female was associated with a 1.58-fold lower risk of developing Relapsing-Remitting Multiple Sclerosis (Wald-test,  $p=0.0036$ ,  $sd = 0.54$ ). In the original dataset, gene expression was assessed in CD4+, CD8+ T cells, and monocytes. To avoid potential confounding effect from paired data (i.e. different measures on the same individual), we focused solely on CD4+ T cell transcriptome to examine the effect of gene expression in the path between gender and MS subtypes. Our analytical model incorporated  $X$  as a binary exposure variable representing gender (0 = male, 1 = female),  $M$  encompassing 18,010 coding gene expressions, and  $Y$  denoting MS subtypes (0 = CIS, 1 = RR).

We conducted a mediation analysis on gene expression to assess the indirect effect of gender on MS subtypes. Due to the small cohort size, we did not identify any significant mediators when applying an FDR control strategy. Alternatively, we opted to perform Gene Set Enrichment Analysis (GSEA) (Korotkevich et al., 2021) to detect biological pathways enriched in our mediation analysis (see Fig 3B). We ranked the mediators using the max-squared test  $p$ -values and screened the mSigDB Hallmark collection of gene sets (Liberzon et al., 2015).



**Figure 3. Effect of gender on MS subtypes.** **A.** Barplot of subtype occurrence according to gender. Total number of individual : 104. **B.** The pathways represent the top 10 selected gene sets. Gene ranks were determined using the  $-\log_{10}$  of the max-squared p-values. NES (Normalized Enrichment Scores) were computed using the GSEA function. Additionally, p-values (pval) and adjusted p-values (padj) via Benjamini-Hochberg correction were associated with each pathway NES.



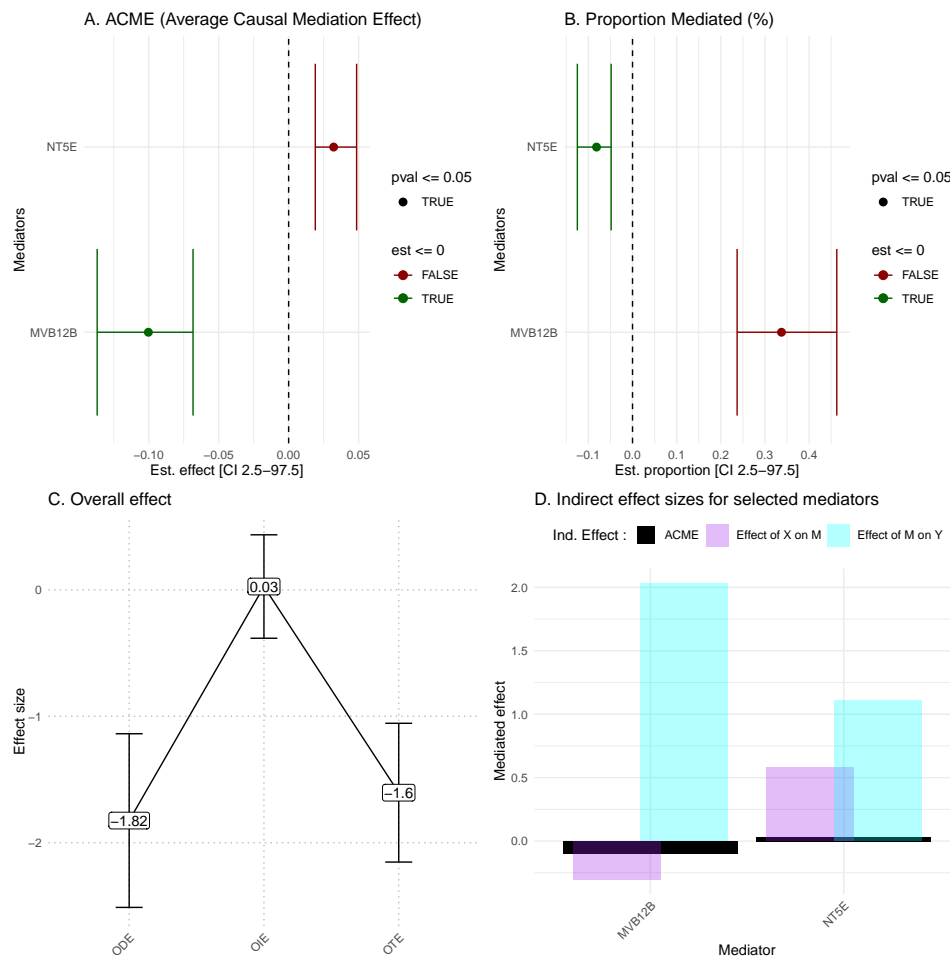
In the Fig 3B, we present the top 10 upregulated pathways identified. It was expected to find hallmarks associated with blood in PBMCs, such as coagulation or heme metabolism. Given the inflammatory nature of MS, it was also reassuring to find immune-related gene sets, such as complement activity or allograft rejection, which correlates with cytolytic activity. Interestingly, we also identified several pathways previously associated with MS disease that thus appear promising for elucidating the causal relationship between gender, gene expression, and MS subtypes : (i) several studies previously established a link between hypoxia and MS disease progression (Halder and Milner, 2020); (ii) some studies have demonstrated the pivotal role of the interleukin 2 receptor pathway in autoimmune response, particularly in MS progression (Peerlings et al., 2021), and (iii) p53 has been associated with immune regulation (Muñoz-Fontela et al., 2016). All of these findings present intriguing avenues for further investigation into the gender impact on MS subtypes.

To further investigate the role of top mediators, despite the low level of significance in the max-squared test, we focused on the top two mediators: *NT5E* (HDMAX2 max-squared p-value =  $2.10 \cdot 10^{-6}$ ) and *MVB12B* (HDMAX2 max-squared p-value =  $5.10 \cdot 10^{-6}$ ). Interestingly, these mediators exhibit opposite effects on the pathway between gender and MS subtypes (Fig 4). *NT5E* is an ecto-5'-nucleotidase known to play a role in immune response deregulation in MS. It has been reported that ectonucleotidases are associated with RR MS in relapsing patients (Álvarez-Sánchez et al., 2019). Notably, *NT5E* expression appears to be elevated in women, thereby increasing the risk of developing RR MS in this population (ACME = +0.03 with a CI of [0.01,0.04]; PM = 8% with a CI of [4,12]). On the other hand, *MVB12B* is a member of the ESCRT-1 complex. This protein was recently identified in a genome-wide protein quantitative trait locus study aimed at identifying drivers of immune-related diseases (The SCALLOP consortium et al., 2023). In our study, *MVB12B* was associated with a reduced risk of RR MS in women (ACME = -0.1 with a CI of [-0.13,-0.07]; PM = 33% with a CI of [23,46]). These findings underscore the ability of high-dimensional analysis to provide insightful understanding into the biological mechanisms underlying observed statistical associations, even with a small cohort size.

## Discussion

In this article, we introduced the *hdmax2* R package, dedicated to high-dimensional mediation analysis. The HDMAX2 method includes unobserved latent factors through a latent factor mixed model approach. It supports the use of exposures of various types and the consideration of both continuous and binary outcomes. We provide a detailed usage vignette and showcase the relevance of our approach through transcriptome and methylome dataset analyses use-cases. The package is available on GitHub for easy access and contribution.

**Figure 4. Summary of the multiple sclerosis use case. A** Estimates of indirect effect (ACME) and **A** proportions of mediated effect (PM) for the top2 mediators. The effect estimate is represented by a dot and its 95% CI by the bar. Symbols correspond to the significance cut off of 5% (square for  $p\text{-value} \geq 0.05$ , circle  $p\text{-value} < 0.05$ ). Colors correspond to the sign of the effect (green for estimated effect  $\leq 0$ , red for estimated effect  $> 0$ ). **C** Effect sizes of Overall Direct Effect (ODE), Overall Indirect Effect (OIE) and Overall Total Effect (OTE). Error bars correspond to standard deviation (ODE and OTE) or confidence interval (OIE). **D** Indirect effect sizes for the selected mediators. Black corresponds to the ACME, violet to the effect of exposure  $X$  on mediator  $M$  in the model  $X \sim M$ , and blue corresponds to the effect of mediator  $M$  on outcome  $Y$  in the model  $Y \sim M + X$ .



While the HDMAX2 method has been optimized for high-dimension analysis, it can naturally be applied outside of this framework, as long as the number of mediators considered is superior to  $K$  (the number of confounders estimated by the latent factor mixed models). The use cases we propose have been selected for demonstration purposes, taking advantage of publicly available data. The choice of mediators to consider should be made by the user based on their specific needs and scientific questions. We provide helper functions in a separate vignette, allowing for conducting mediator selection under FDR control, or for the aggregation of p-values to study regions of interest. Aggregated Methylated Regions (AMR) are a typical example of regions of interest, when studying DNA methylation mediated effect. AMR are made of CpG with significant p-value, in a given genomic region, they can be viewed as the parallel of Differentially Methylated Regions in classical EWAS. Although mediation analyses can establish a statistical association between an exposure, a mediator and an outcome, they do not guarantee a causal role in the biological processes observed. The user must then be careful not to over-interpret his results, and to build his models taking into account the sequentiality of the elements observed and the interactions between the different variables included in the model.

We believe that the HDMAX2 method and its implementation as a package could find wide applications in the fields of environmental and clinical epidemiology, as well as in computational biology approaches aimed at gaining deeper insights into the biological background of diseases such as cancer, for which omics datasets are readily available in the public domain. It is worth noting that recent studies have proposed methods for conducting mediation analysis in the context of survival outcomes (Clark-Boucher et al., 2023; Luo et al., 2020; H Zhang, Zheng, Hou, et al., 2021). This is a significant demand in healthcare, and we are considering including this capability in the `hdmax2` package once we have conducted a comprehensive analysis of the methodological approach to implement.

## Acknowledgements

We thank Florent Chuffart for his critical reading of the manuscript. We are grateful to the discussion group *Meth*<sup>2</sup>, Lucile Broseus, Ariane Guilbert and Claire-Cécile Barrot for their helpful discussions.

## Fundings

This work was funded by the project THEMA, from "Appel à projets IRGA 2021-2022" from the University Grenoble Alpes and by the French Agency for National Research (ETAPE // ANR-18-CE36-0005 and CauseHet // ANR-22-CE45-0030). It has also been carried out with financial support from ITMO Cancer of Aviesan within the framework of the 2021-2030 Cancer Control Strategy, on funds administered by Inserm (ACACIA project AAP-MIC-2021).

## Conflict of interest disclosure

The authors declare that they comply with the PCI rule of having no financial conflicts of interest in relation to the content of the article.

## Data, script, code, and supplementary information availability

Data are available online on public repositories. Script and codes are available online: <https://github.com/bcm-uga/hdmax2>.

## References

- Álvarez-Sánchez N, I Cruz-Chamorro, M Díaz-Sánchez, PJ Lardone, JM Guerrero, and A Carrillo-Vico (Feb. 2019). Peripheral CD39-expressing T regulatory cells are increased and associated with relapsing-remitting multiple sclerosis in relapsing patients. en. *Scientific Reports* 9. Publisher: Nature Publishing Group, 2302. ISSN: 2045-2322. <https://doi.org/10.1038/s41598-019-38897-w>.
- Baron RM and DA Kenny (Dec. 1986). The moderator-mediator variable distinction in social psychological research: conceptual, strategic, and statistical considerations. eng. *Journal of Personality and Social Psychology* 51, 1173–1182. ISSN: 0022-3514. <https://doi.org/10.1037//0022-3514.51.6.1173>.
- Blum MGB, L Valeri, O François, S Cadiou, V Siroux, J Lepeule, and R Slama (May 2020). Challenges Raised by Mediation Analysis in a High-Dimension Setting. eng. *Environmental Health Perspectives* 128, 55001. ISSN: 1552-9924. <https://doi.org/10.1289/EHP6240>.

Caye K, B Jumentier, J Lepeule, and O François (Apr. 2019). LFMM 2: Fast and Accurate Inference of Gene-Environment Associations in Genome-Wide Studies. *Molecular Biology and Evolution* 36, 852–860. ISSN: 0737-4038. <https://doi.org/10.1093/molbev/msz008>.

Clark-Boucher D, X Zhou, J Du, Y Liu, BL Needham, JA Smith, and B Mukherjee (Feb. 2023). Methods for Mediation Analysis with High-Dimensional DNA Methylation Data: Possible Choices and Comparison. en. Pages: 2023.02.10.23285764. <https://doi.org/10.1101/2023.02.10.23285764>.

Dai JY, JL Stanford, and M LeBlanc (Jan. 2022). A Multiple-Testing Procedure for High-Dimensional Mediation Hypotheses. *Journal of the American Statistical Association* 117. Publisher: Taylor & Francis \_eprint: <https://doi.org/10.1080/01621459.2020.1765785>, 198–213. ISSN: 0162-1459. <https://doi.org/10.1080/01621459.2020.1765785>.

Djordjilović V, J Hemerik, and M Thoresen (2022). On optimal two-stage testing of multiple mediators. en. *Biometrical Journal* 64. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/bimj.202100190>, 1090–1108. ISSN: 1521-4036. <https://doi.org/10.1002/bimj.202100190>.

Dolgalev I (2024). en.

Halder SK and R Milner (Dec. 2020). Hypoxia in multiple sclerosis; is it the chicken or the egg? *Brain* 144, 402–410. ISSN: 0006-8950. <https://doi.org/10.1093/brain/awaa427>.

Harbo HF, R Gold, and M Tintoré (July 2013). Sex and gender issues in multiple sclerosis. *Therapeutic Advances in Neurological Disorders* 6, 237–248. ISSN: 1756-2856. <https://doi.org/10.1177/1756285613488434>.

Imai K, L Keele, and D Tingley (2010). A general approach to causal mediation analysis. en. *Psychological Methods* 15, 309–334. ISSN: 1939-1463, 1082-989X. <https://doi.org/10.1037/a0020761>.

Jumentier B, CC Barrot, M Estavoyer, J Tost, B Heude, O François, and J Lepeule (2023). High-Dimensional Mediation Analysis: A New Method Applied to Maternal Smoking, Placental DNA Methylation, and Birth Outcomes. *Environmental Health Perspectives* 131 (). Publisher: Environmental Health Perspectives, 047011. <https://doi.org/10.1289/EHP11559>.

Kim K, AK Pröbstel, R Baumann, J Dyckow, J Landefeld, E Kogl, L Madireddy, R Loudermilk, EL Eggers, S Singh, SJ Caillier, SL Hauser, BAC Cree, UCSF MS-EPIC Team, L Schirmer, MR Wilson, and SE Baranzini (Feb. 2021). Cell type-specific transcriptomics identifies neddylation as a novel therapeutic target in multiple sclerosis. *Brain* 144, 450–461. ISSN: 0006-8950. <https://doi.org/10.1093/brain/awaa421>.

Korotkevich G, V Sukhov, N Budin, B Shpak, MN Artyomov, and A Sergushichev (Feb. 2021). Fast gene set enrichment analysis. en. Pages: 060012 Section: New Results. <https://doi.org/10.1101/060012>.

Liberzon A, C Birger, H Thorvaldsdóttir, M Ghandi, JP Mesirov, and P Tamayo (Dec. 2015). The Molecular Signatures Database (MSigDB) hallmark gene set collection. *Cell systems* 1, 417–425. ISSN: 2405-4712. <https://doi.org/10.1016/j.cels.2015.12.004>.

Love MI, W Huber, and S Anders (Dec. 2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. en. *Genome Biology* 15, 550. ISSN: 1474-760X. <https://doi.org/10.1186/s13059-014-0550-8>.

Luo C, B Fa, Y Yan, Y Wang, Y Zhou, Y Zhang, and Z Yu (Apr. 2020). High-dimensional mediation analysis in survival models. en. *PLOS Computational Biology* 16. Ed. by Althouse B, e1007768. ISSN: 1553-7358. <https://doi.org/10.1371/journal.pcbi.1007768>.

Muñoz-Fontela C, A Mandinova, SA Aaronson, and SW Lee (Dec. 2016). Emerging roles of p53 and other tumour-suppressor genes in immune regulation. en. *Nature Reviews Immunology* 16. Publisher: Nature Publishing Group, 741–750. ISSN: 1474-1741. <https://doi.org/10.1038/nri.2016.99>.

Orrantia-Borunda E, P Anchondo-Nuñez, LE Acuña-Aguilar, FO Gómez-Valles, and CA Ramírez-Valdespino (2022). Subtypes of Breast Cancer. eng. In: *Breast Cancer*. Ed. by Mayrovitz HN. Brisbane (AU): Exon Publications. ISBN: 978-0-645-33203-2.

Peerlings D, M Mimpen, and J Damoiseaux (2021). The IL-2 - IL-2 receptor pathway: Key to understanding multiple sclerosis. eng. *Journal of Translational Autoimmunity* 4, 100123. ISSN: 2589-9090. <https://doi.org/10.1016/j.jtauto.2021.100123>.

Richiardi L, R Bellocco, and D Zugna (Oct. 2013). Mediation analysis in epidemiology: methods, interpretation and bias. *International Journal of Epidemiology* 42, 1511–1519. ISSN: 0300-5771. <https://doi.org/10.1093/ije/dyt127>.

Sampson JN, SM Boca, SC Moore, and R Heller (July 2018). FWER and FDR control when testing multiple mediators. *Bioinformatics* 34, 2418–2424. ISSN: 1367-4803. <https://doi.org/10.1093/bioinformatics/bty064>.

Sobel ME (1982). Asymptotic Confidence Intervals for Indirect Effects in Structural Equation Models. *Sociological Methodology* 13. Publisher: [American Sociological Association, Wiley, Sage Publications, Inc.], 290–312. ISSN: 0081-1750. <https://doi.org/10.2307/270723>.

The SCALLOP consortium, JH Zhao, D Stacey, N Eriksson, E Macdonald-Dunlop, ÅK Hedman, A Kalnapenkis, S Enroth, D Cozzetto, J Digby-Bell, J Marten, L Folkersen, C Herder, L Jonsson, SE Bergen, C Geiger, EJ Needham, P Surendran, Estonian Biobank Research Team, DS Paul, O Polasek, B Thorand, H Grallert, M Roden, U Võsa, T Esko, C Hayward, Å Johansson, U Gyllensten, N Powell, O Hansson, N Mattsson-Carlsson, PK Joshi, J Danesh, L Padyukov, L Klareskog, M Landén, JF Wilson, A Siegbahn, L Wallentin, A Mälarstig, AS Butterworth, and JE Peters (Mar. 2023). *Mapping pQTLs of circulating inflammatory proteins identifies drivers of immune-related disease risk and novel therapeutic targets*. en. preprint. Genetic and Genomic Medicine. <https://doi.org/10.1101/2023.03.24.23287680>.

Tingley D, T Yamamoto, K Hirose, L Keele, and K Imai (2014). **mediation** : R Package for Causal Mediation Analysis. en. *Journal of Statistical Software* 59. ISSN: 1548-7660. <https://doi.org/10.18637/jss.v059.i05>.

Voskuhl RR, K Patel, F Paul, SM Gold, M Scheel, J Kuchling, G Cooper, S Assemer, C Chien, AU Brandt, CE Meyer, and A MacKenzie-Graham (Aug. 2020). Sex differences in brain atrophy in multiple sclerosis. *Biology of Sex Differences* 11, 49. ISSN: 2042-6410. <https://doi.org/10.1186/s13293-020-00326-3>.

Yin Q, H Ma, Y Dong, S Zhang, J Wang, J Liang, L Mao, L Zeng, X Xiong, X Chen, J Wang, and X Zheng (Dec. 2024). The integration of multidisciplinary approaches revealed PTGES3 as a novel drug target for breast cancer treatment. en. *Journal of Translational Medicine* 22. Number: 1 Publisher: BioMed Central, 1–16. ISSN: 1479-5876. <https://doi.org/10.1186/s12967-024-04899-0>.

Zeng P, Z Shao, and X Zhou (May 2021). Statistical methods for mediation analysis in the era of high-throughput genomics: Current successes and future challenges. *Computational and Structural Biotechnology Journal* 19, 3209–3224. ISSN: 2001-0370. <https://doi.org/10.1016/j.csbj.2021.05.042>.

Zhang H, Y Zheng, L Hou, C Zheng, and L Liu (Nov. 2021). Mediation analysis for survival data with high-dimensional mediators. eng. *Bioinformatics (Oxford, England)* 37, 3815–3821. ISSN: 1367-4811. <https://doi.org/10.1093/bioinformatics/btab564>.

Zhang H, Y Zheng, Z Zhang, T Gao, B Joyce, G Yoon, W Zhang, J Schwartz, A Just, E Colicino, P Vokonas, L Zhao, J Lv, A Baccarelli, L Hou, and L Liu (Oct. 2016). Estimating and testing high-dimensional mediation effects in epigenetic studies. *Bioinformatics* 32, 3150–3154. ISSN: 1367-4803. <https://doi.org/10.1093/bioinformatics/btw351>.