

# Supplementary Materials

## Machine learning-based analysis of glioma grades reveals co-enrichment

Mateusz Garbulowski<sup>1,2,\*</sup>, Karolina Smolinska<sup>1</sup>, Uğur Çabuk<sup>1,3,4</sup>, Sara A. Yones<sup>1</sup>, Ludovica Celli<sup>1,5,6</sup>, Esma Nur Yaz<sup>1,7</sup>, Fredrik Barrenäs<sup>1,8</sup>, Klev Diamanti<sup>1,9</sup>, Claes Wadelius<sup>9</sup> and Jan Komorowski<sup>1,8,10,11,\*</sup>

<sup>1</sup>Department of Cell and Molecular Biology, Uppsala University, Uppsala, Sweden

<sup>2</sup>Department of Biochemistry and Biophysics, Science for Life Laboratory, Stockholm University, Solna, Sweden

<sup>3</sup>Polar Terrestrial Environmental Systems, Alfred Wegener Institute Helmholtz Centre for Polar and Marine Research, Potsdam, Germany

<sup>4</sup>Institute of Biochemistry and Biology, University of Potsdam, Potsdam, Germany

<sup>5</sup>Institute of Molecular Genetics Luigi Luca Cavalli-Sforza, National Research Council, Pavia, Italy

<sup>6</sup>Department of Biology and Biotechnology, University of Pavia, Pavia Italy

<sup>7</sup>Department of Biomedical Engineering and Bioinformatics, The Graduate School of Engineering and Natural Sciences, Istanbul Medipol University, Istanbul, Turkey

<sup>8</sup>Washington National Primate Research Center, Seattle, WA, USA

<sup>9</sup>Department of Immunology, Genetics and Pathology, Uppsala University, Uppsala, Sweden

<sup>10</sup>Swedish Collegium for Advanced Study, Uppsala, Uppsala, Sweden

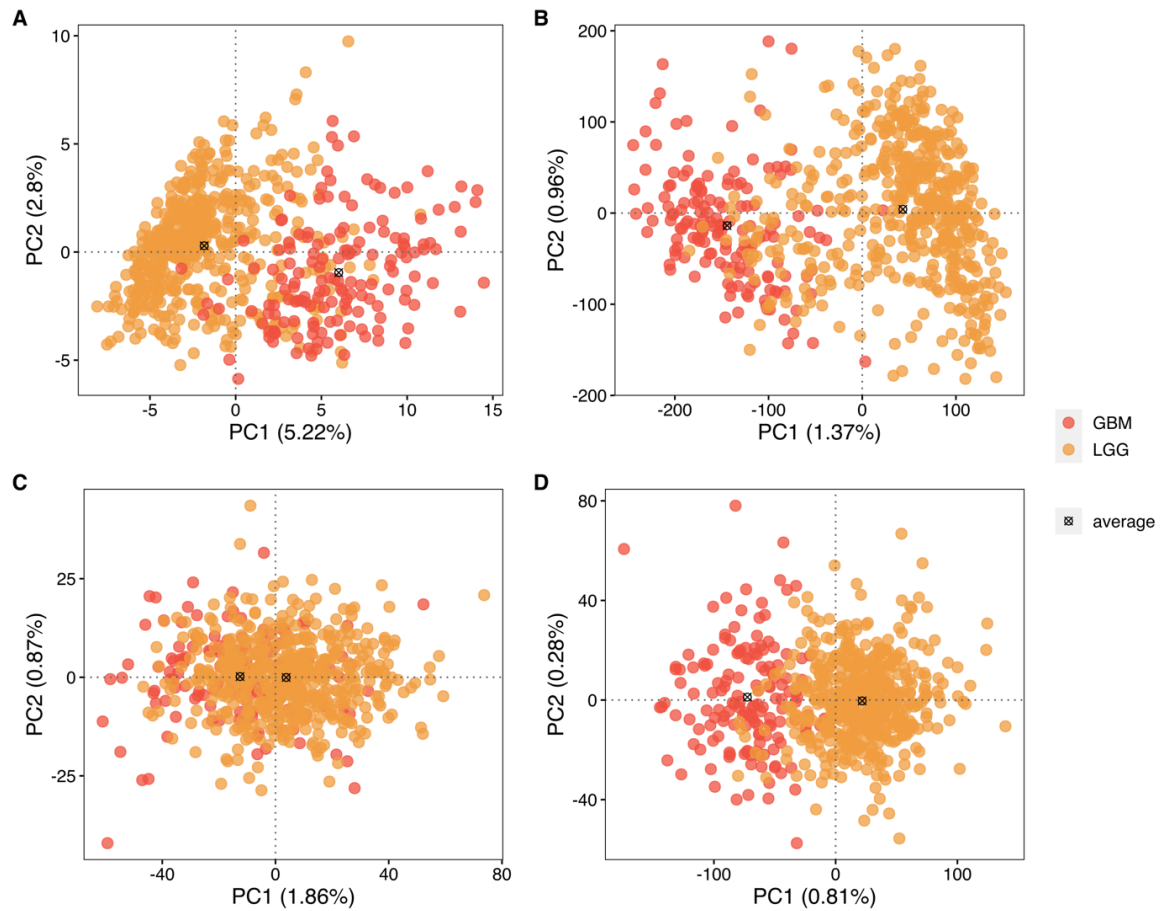
<sup>11</sup>Institute of Computer Science, Polish Academy of Sciences, Warsaw, Poland

\*Correspondence: M.G. [mateuszgarbulowski@gmail.com](mailto:mateuszgarbulowski@gmail.com); J.K. [jan.komorowski@icm.uu.se](mailto:jan.komorowski@icm.uu.se)

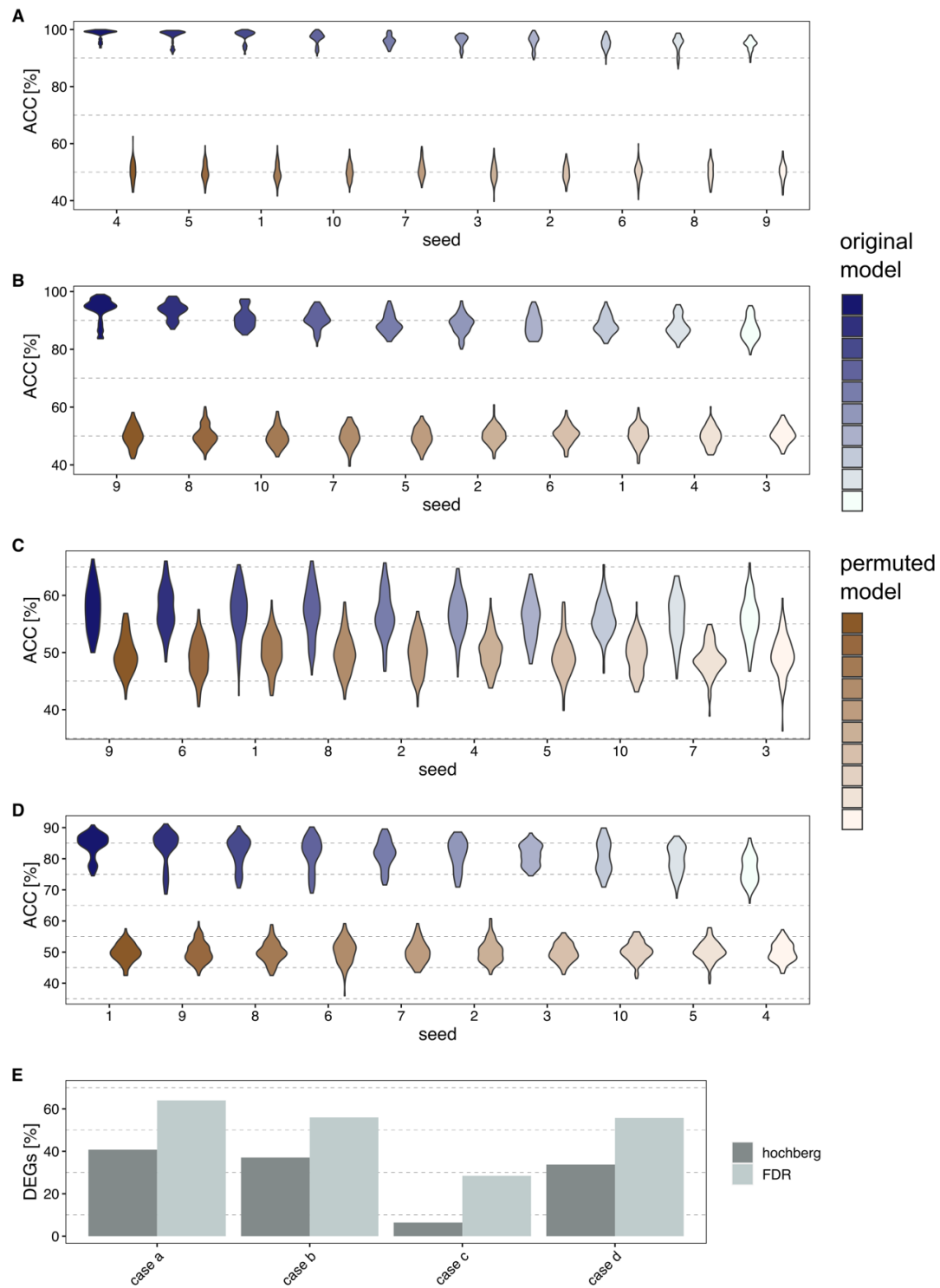
## Table of Contents

<i>Supplementary Figures .....</i>	<b>3</b>
<i>Supplementary Captions for Tables S1-S8.....</i>	<b>25</b>
<i>Supplementary References for Table S1.....</i>	<b>26</b>

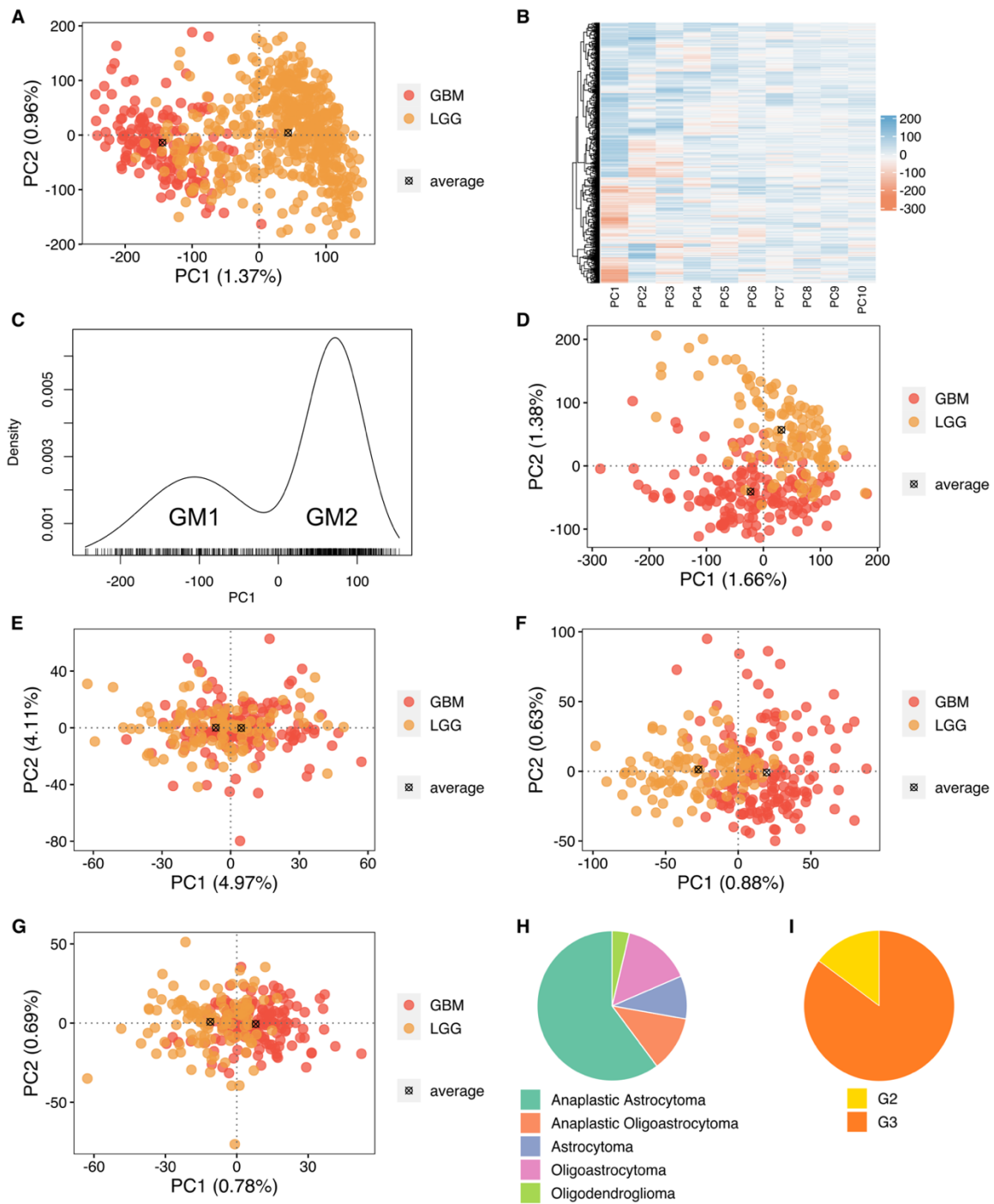
## Supplementary Figures



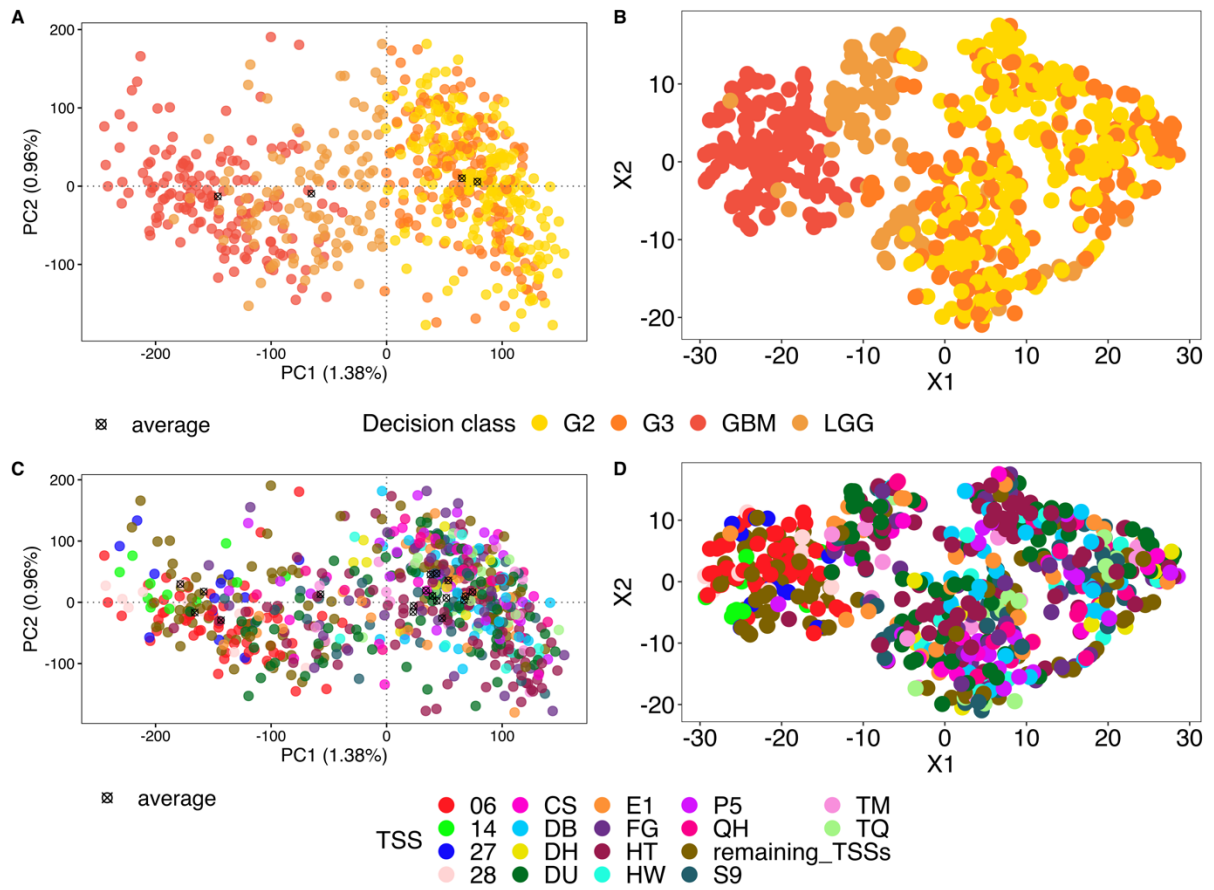
**Figure S1.** Principal component analysis (PCA) to investigate the variance for The Cancer Genome Atlas (TCGA) cohorts for Lower Grade Glioma (LGG) and Glioblastoma Multiforme (GBM). **(A)** Original data from the GDC repository. **(B)** Unified cohorts from the UCSC Xena repository. **(C)** Batch effect correction of unified cohorts for the *leek* parameter with 532 surrogate variables. **(D)** Batch effect correction of unified cohorts for the *be* parameter with 48 surrogate variables.



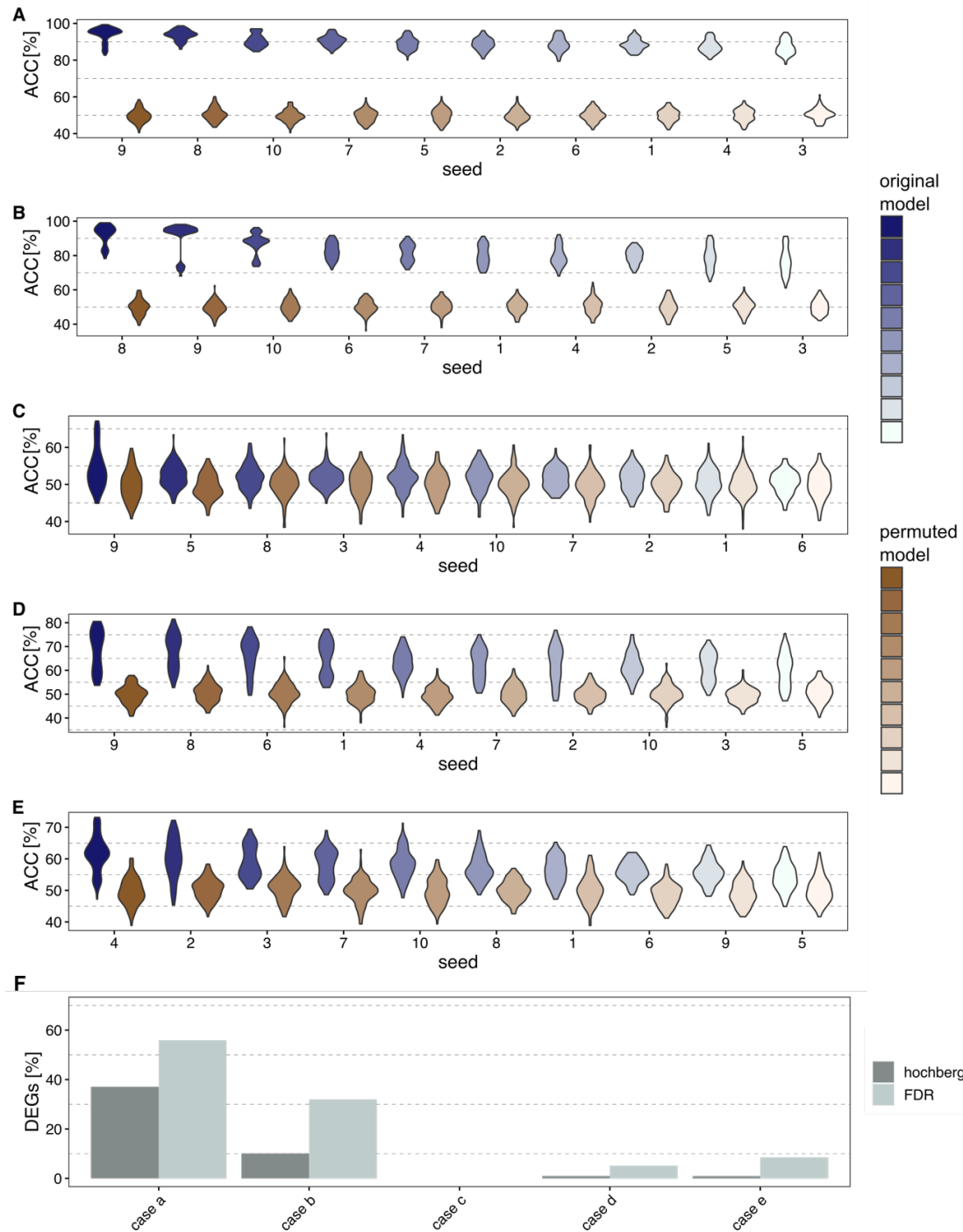
**Figure S2.** Machine learning analysis was performed on 100 randomly selected genes with 10 various seed values to investigate the bias for TCGA cohorts of LGGs and GBMs. All five machine learning methods were applied, viz. SMO, IBk, Bagging, J48 and JRip (see section ML evaluation). (A) Original data from the GDC repository. (B) Unified cohorts from the UCSC Xena repository. (C) Batch effect correction of unified cohorts for the *leek* parameter with 532 surrogate variables. (D) Batch effect correction of unified cohorts for the *be* parameter with 48 surrogate variables. (E) Student's *t*-test results for particular cases a-d corresponding to subplots (A-D).



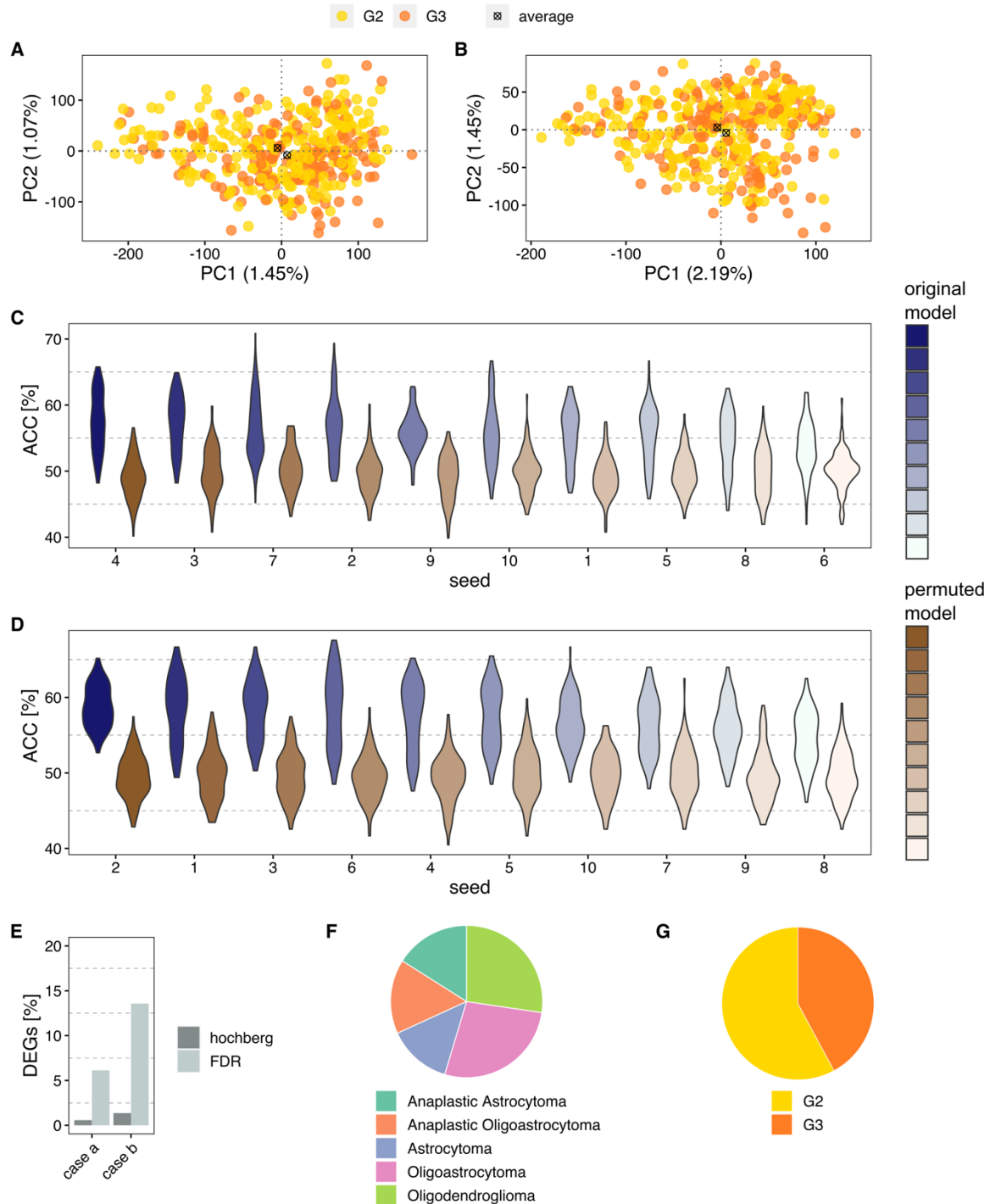
**Figure S3.** Data evaluation to localize and remove bias from the glioma cohorts (A) PCA for unified cohorts from the UCSC Xena repository. (B) Clustering of 1-10 PCs of unified cohorts from the UCSC Xena repository. (C) Gaussian mixtures (GM) were detected from PC1. In GM2-based sample sets, two GBMs were included in the mixture and excluded in further analyses. (D) PCA was performed on LGG and GBM samples selected based on GM1. (E) Batch effect correction performed on LGG and GBM samples selected based on GM1 for the *leek* parameter with 232 surrogate variables. (F) Batch effect correction for the *be* parameter with 30 surrogate variables performed on samples selected based on GM1. (G) The final data set was based on GM1 after batch effect correction and filtration for protein-coding genes (H-I) Evaluation of LGG samples included in GM1 based data set.



**Figure S4.** The global and local analysis of the data structure for TCGA glioma cohorts. **(A)** PCA and **(B)** t-SNE for decision class selected based on GMs. **(C)** PCA and **(D)** t-SNE for tissue source sites (TSSs). t-SNE plots were created for seed 1. Remaining TSSs included groups of samples equal to or less than 10.

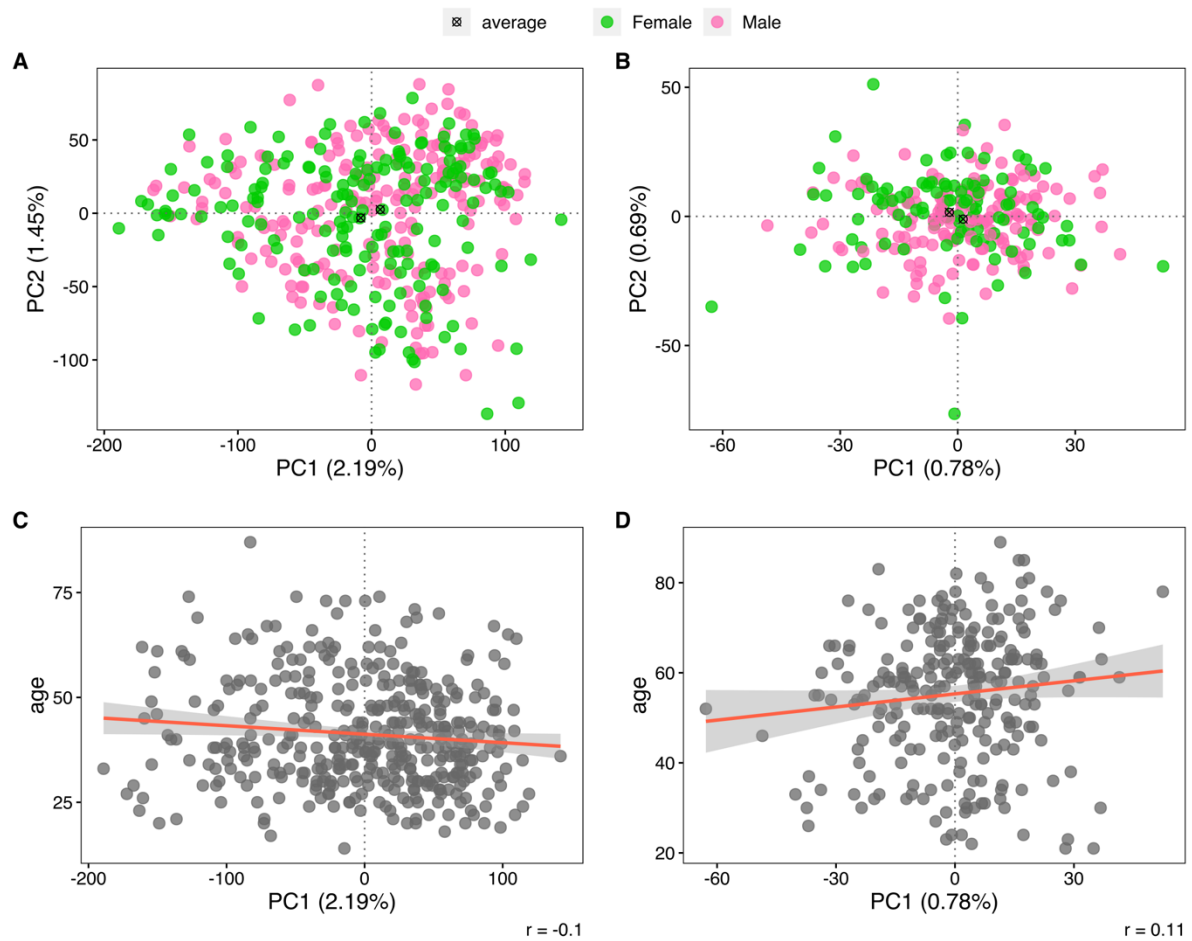


**Figure S5.** Machine learning analysis was performed on 100 randomly selected genes with 10 various seed values to investigate the bias for TCGA cohorts of LGGs and GBMs. All five machine learning methods were applied, viz. SMO, IBk, Bagging, J48 and JRip (see section ML evaluation). (A) unified cohorts from the UCSC Xena repository. (B) GM2 samples from unified cohorts from the UCSC Xena repository. (C) Batch effect correction of 259 GM2-based samples for the *leek* parameter with 232 surrogate variables. (D) Batch effect correction of 259 GM2-based samples for the *be* parameter with 30 surrogate variables. (E) The final data set is based on GM1 after batch effect correction and filtration for protein-coding genes. (F) Student's *t*-test results for particular cases a-e corresponding to subplots (A-E).

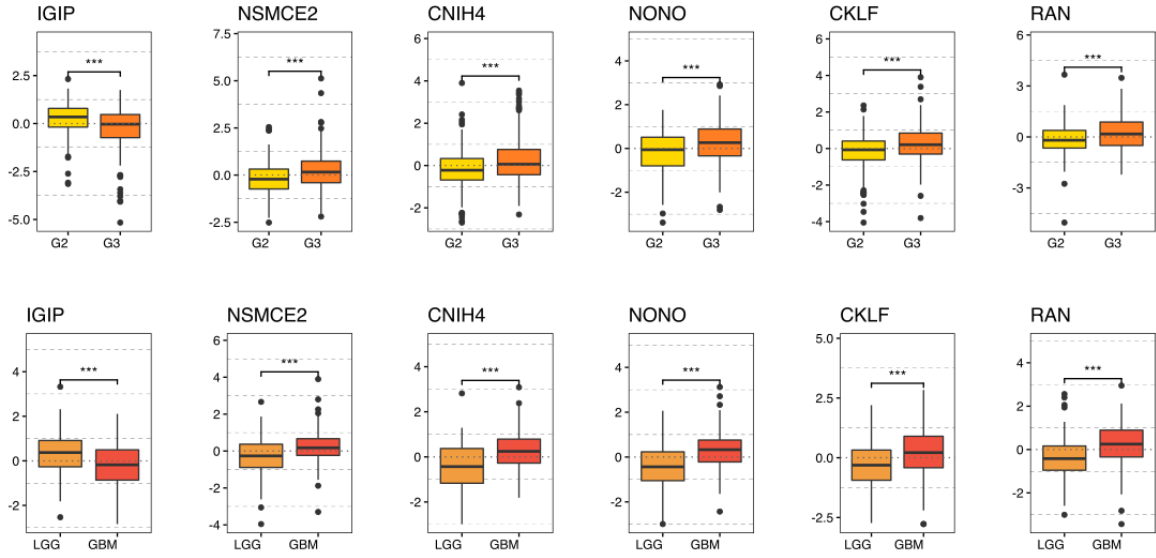


**Figure S6.** The evaluation of samples in GM2 consisting of the vast majority of LGGs. **(A)** PCA for LGGs based on GM2 samples. **(B)** PCA for LGGs based on GM2 samples for protein-coding genes only. **(C)** Machine learning evaluation of GM2 samples for discerning between GII and GIII. **(D)** Machine learning evaluation of GM2 samples for discerning GII and GIII for protein-coding genes only. **(E)** Evaluation of fraction of statistically significant genes discerning between GII and GIII with all genes (case a) and protein-coding genes only (case b). **(F)** The proportion of grades within a GM2-based data set for LGGs. For machine learning evaluation **(C, D)** all five methods were applied, viz. SMO, IBk, Bagging, J48 and JRip (see section ML evaluation).

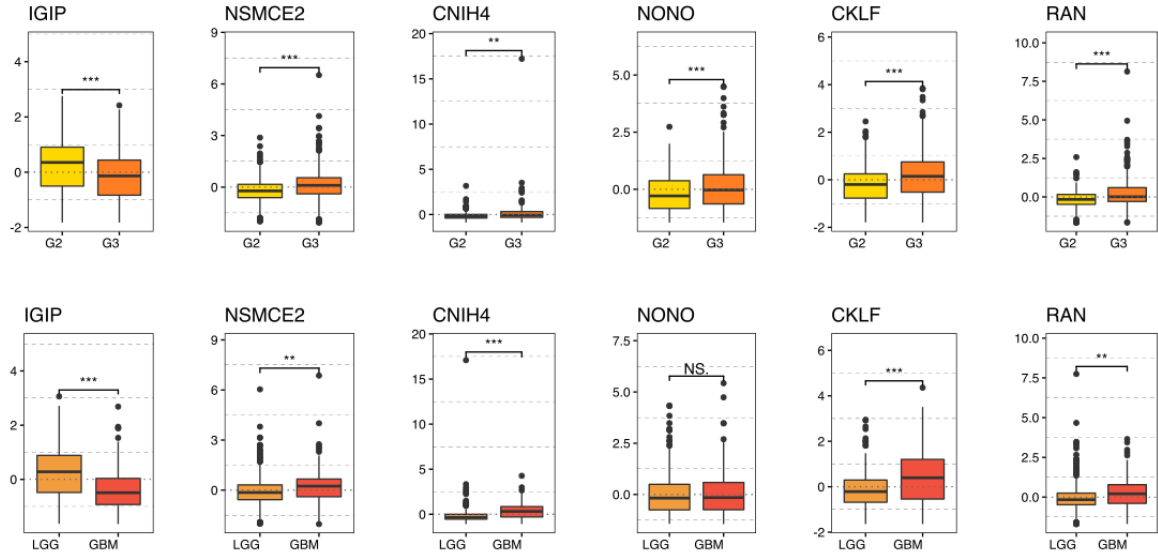




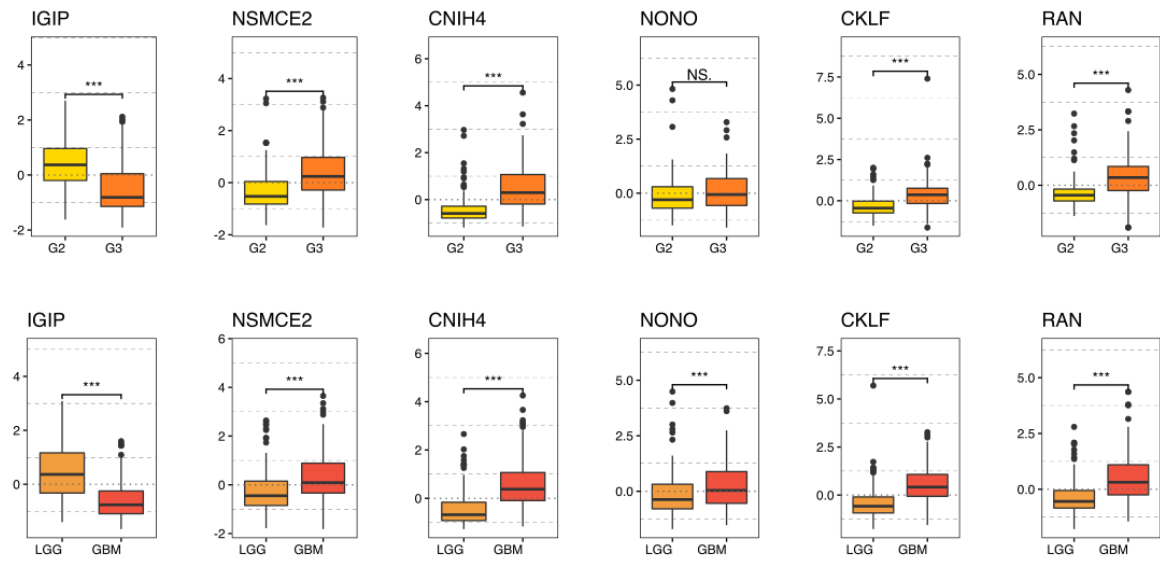
**Figure S7.** The evaluation of other clinical factors within GM1- and GM2-based data sets. **(A)** PCA for GM1 for LGGs and their sex information. **(B)** PCA for GM2 for LGGs and GBMs and their sex information. **(C)** The variation of the age of LGG samples is based on GM1. **(D)** The variation of the age of LGG and GBM samples is based on GM2. Pearson correlation ( $r$ ) value is marked in the plot caption.



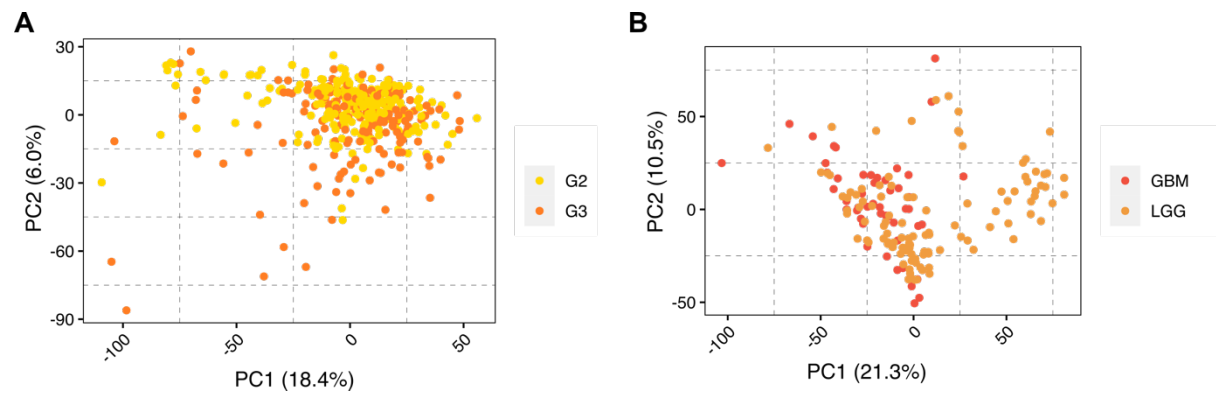
**Figure S8.** Expression profiles for six common differentially expressed genes (DEGs) were selected based on the intersection of the DEGs list between GII vs. GIII and LGG vs. GBM. Gene expression profiles were generated based on TCGA cohorts. *P* values were marked on boxplots as ns ( $P > 0.05$ ), \* ( $P \leq 0.05$ ), \*\* ( $P \leq 0.01$ ), and \*\*\* ( $P \leq 0.001$ ).



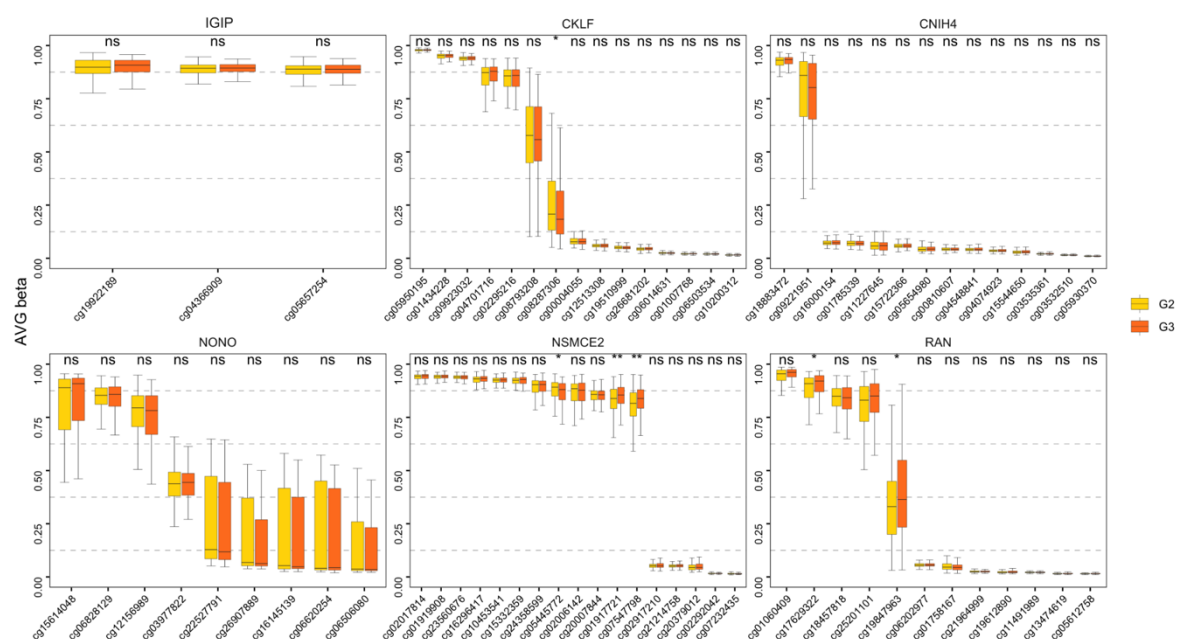
**Figure S9.** Expression profiles for six common DEGs were selected based on the intersection of DEGs list between GII vs. GIII and LGG vs. GBM. Gene expression profiles were generated based on the Chinese Glioma Genome Atlas (CGGA) batch 1 cohort.  $P$  values were marked on boxplots as ns ( $P > 0.05$ ), \* ( $P \leq 0.05$ ), \*\* ( $P \leq 0.01$ ), and \*\*\* ( $P \leq 0.001$ ).



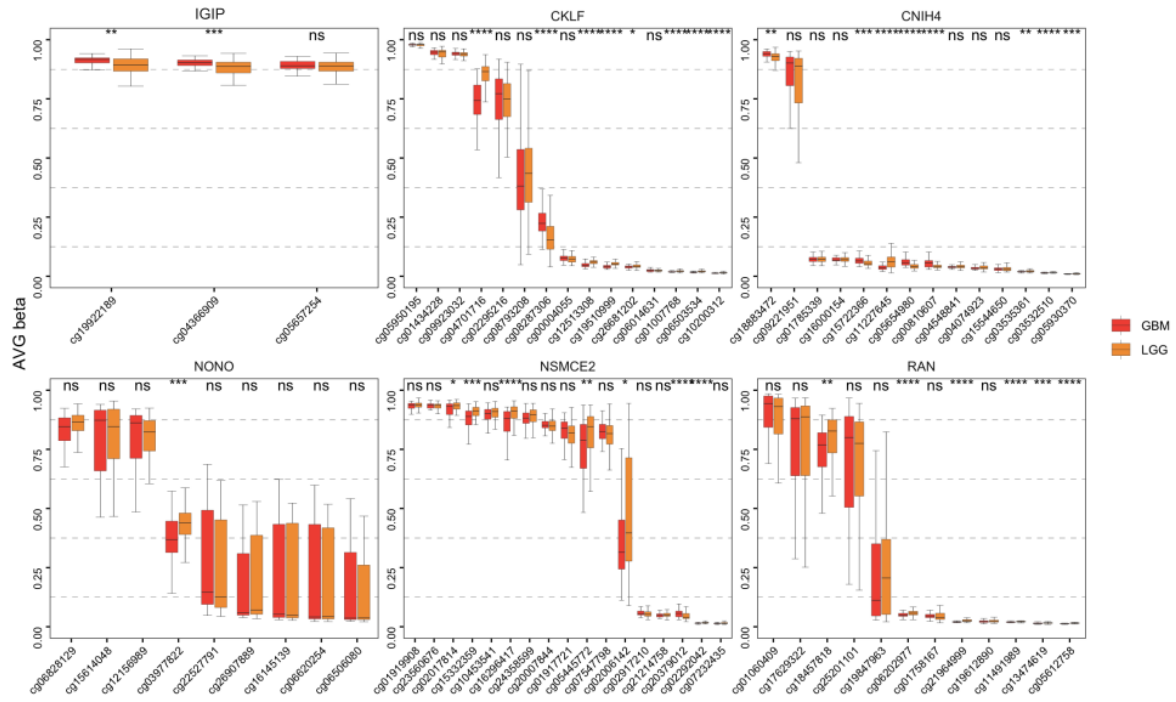
**Figure S10.** Expression profiles for six common DEGs were selected based on the intersection of DEGs list between GII vs. GIII and LGG vs. GBM. Gene expression profiles were generated based on the CGGA batch 2 cohort.  $P$  values were marked on boxplots as ns ( $P > 0.05$ ), \* ( $P \leq 0.05$ ), \*\* ( $P \leq 0.01$ ), and \*\*\* ( $P \leq 0.001$ ).



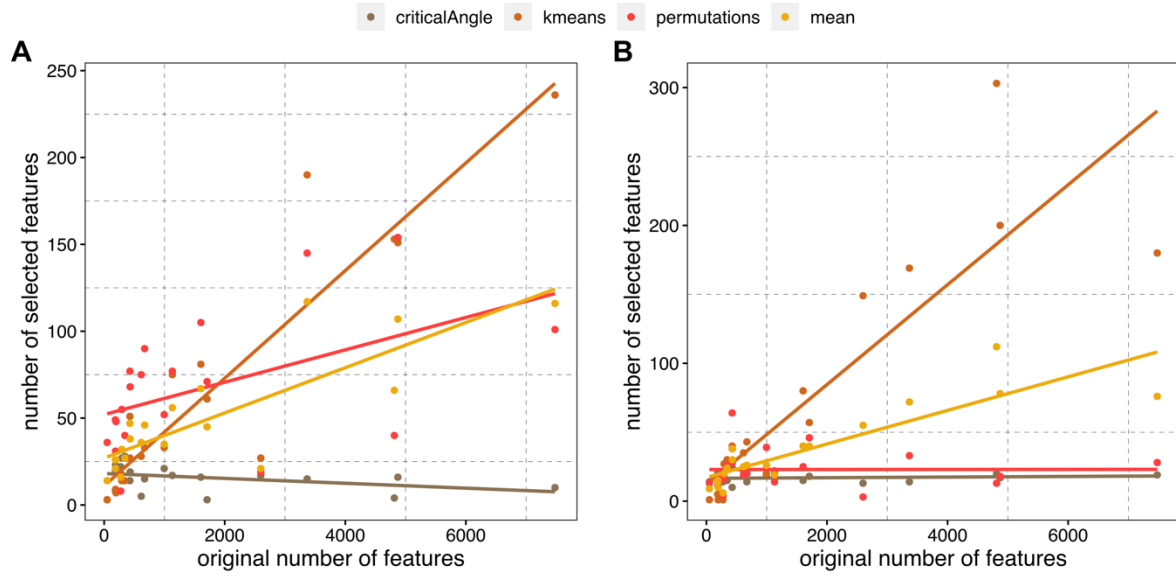
**Figure S11.** Variability in DNA methylation data from TCGA based on samples selected upon GM modeling for **(A)** GM2 and **(B)** GM1. Explained variation is given in parenthesis.



**Figure S12.** DNA methylation status of six common DEGs for GII vs. GIII.  $P$  values were marked on boxplots as ns ( $P > 0.05$ ), \* ( $P \leq 0.05$ ), \*\* ( $P \leq 0.01$ ), and \*\*\* ( $P \leq 0.001$ ).

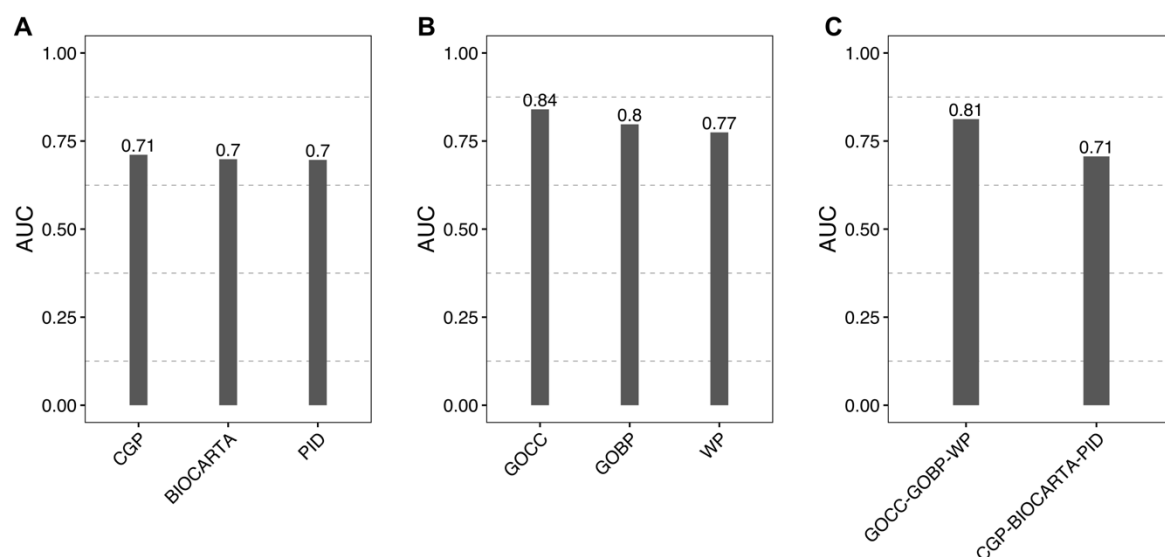


**Figure S13.** DNA methylation status of six common DEGs for LGG vs. GBM.  $P$  values were marked on boxplots as ns ( $P > 0.05$ ), \* ( $P \leq 0.05$ ), \*\* ( $P \leq 0.01$ ), and \*\*\* ( $P \leq 0.001$ ).

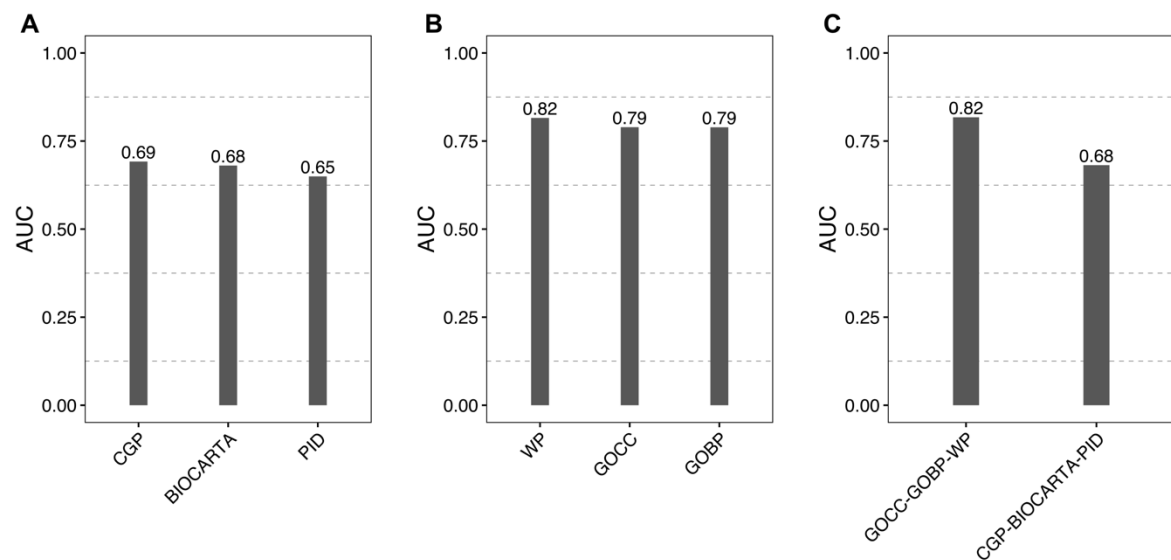


**Figure S14.** The Monte Carlo feature selection (MCFS) evaluation of thresholds was performed for choosing a threshold for selecting top features. **(A)** All 20 MSigDB collections for GII vs. GIII. **(B)** All 20 MSigDB collections for LGG vs. GBM.

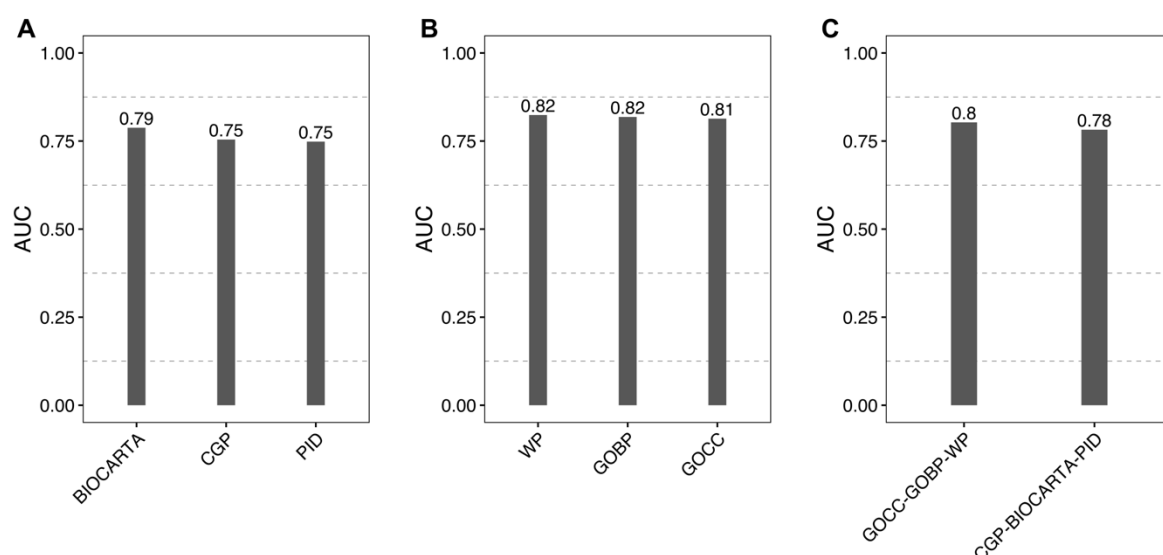




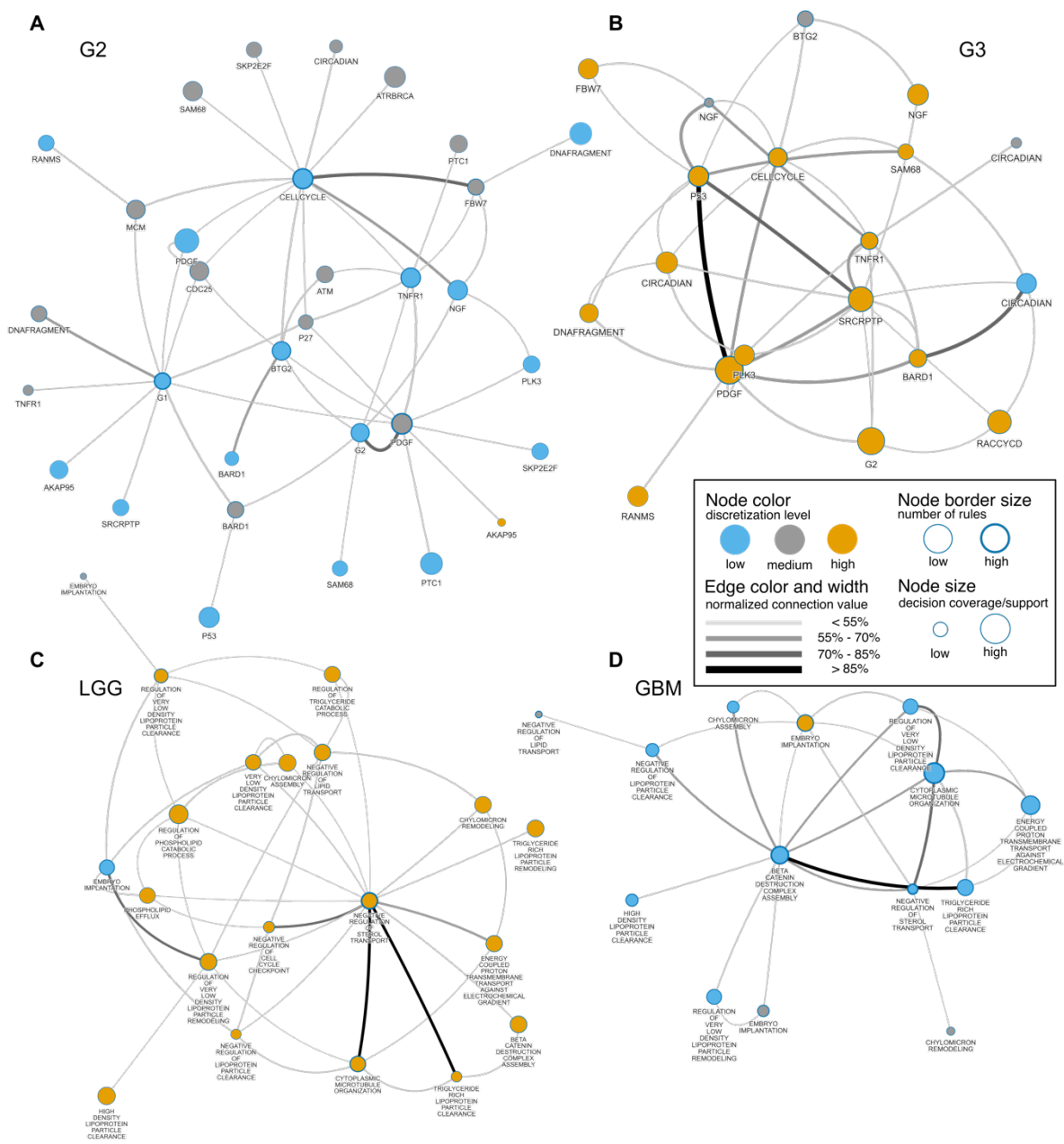
**Figure S15.** Evaluation of MSigDB collections with R.ROSETTA rule-based learning for classifying glioma grades using ssGSEA scores based on TCGA cohorts. **(A)** Top three MSigDB collections for classifying GII vs. GIII. **(B)** Top three MSigDB collections for classifying LGG vs. GBM. **(C)** Merged top three MSigDB collections for two RBMs classifying GII vs. GIII (right bar), and LGG vs. GBM (left bar).



**Figure S16.** Evaluation of MSigDB collections with R.ROSETTA rule-based learning for classifying glioma grades using ssGSEA scores based on CGGA batch 1. **(A)** Top three MSigDB collections for classifying GII vs. GIII. **(B)** Top three MSigDB collections for classifying LGG vs. GBM. **(C)** Merged top three MSigDB collections for two RBMs classifying GII vs. GIII (right bar), and LGG vs. GBM (left bar).

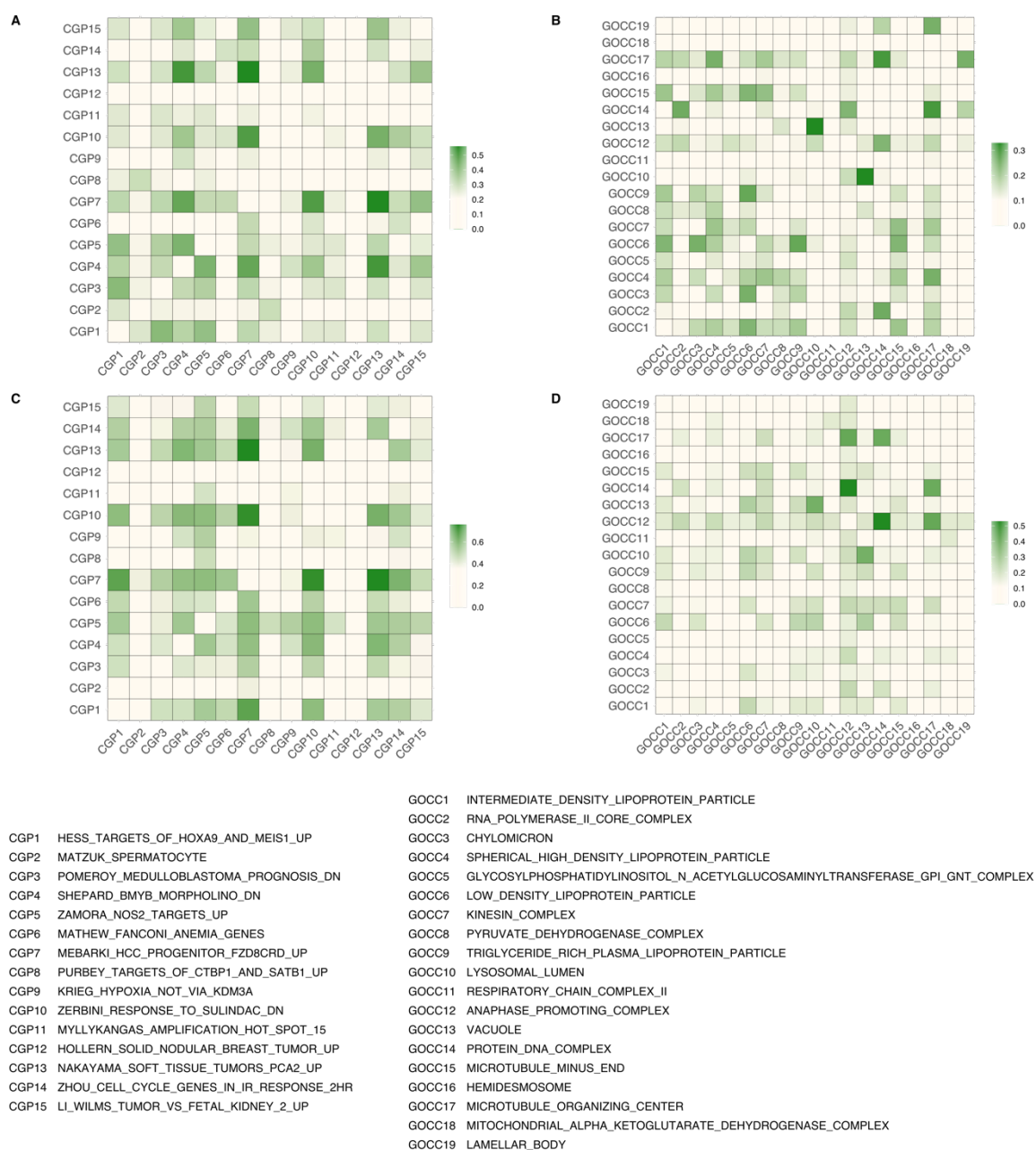


**Figure S17.** Evaluation of MSigDB collections with R.ROSETTA rule-based learning for classifying glioma grades using ssGSEA scores based on CGGA batch 2. **(A)** Top three MSigDB collections for classifying GII vs. GIII. **(B)** Top three MSigDB collections for classifying LGG vs. GBM. **(C)** Merged top three MSigDB collections for two RBMs classifying GII vs. GIII (right bar), and LGG vs. GBM (left bar).

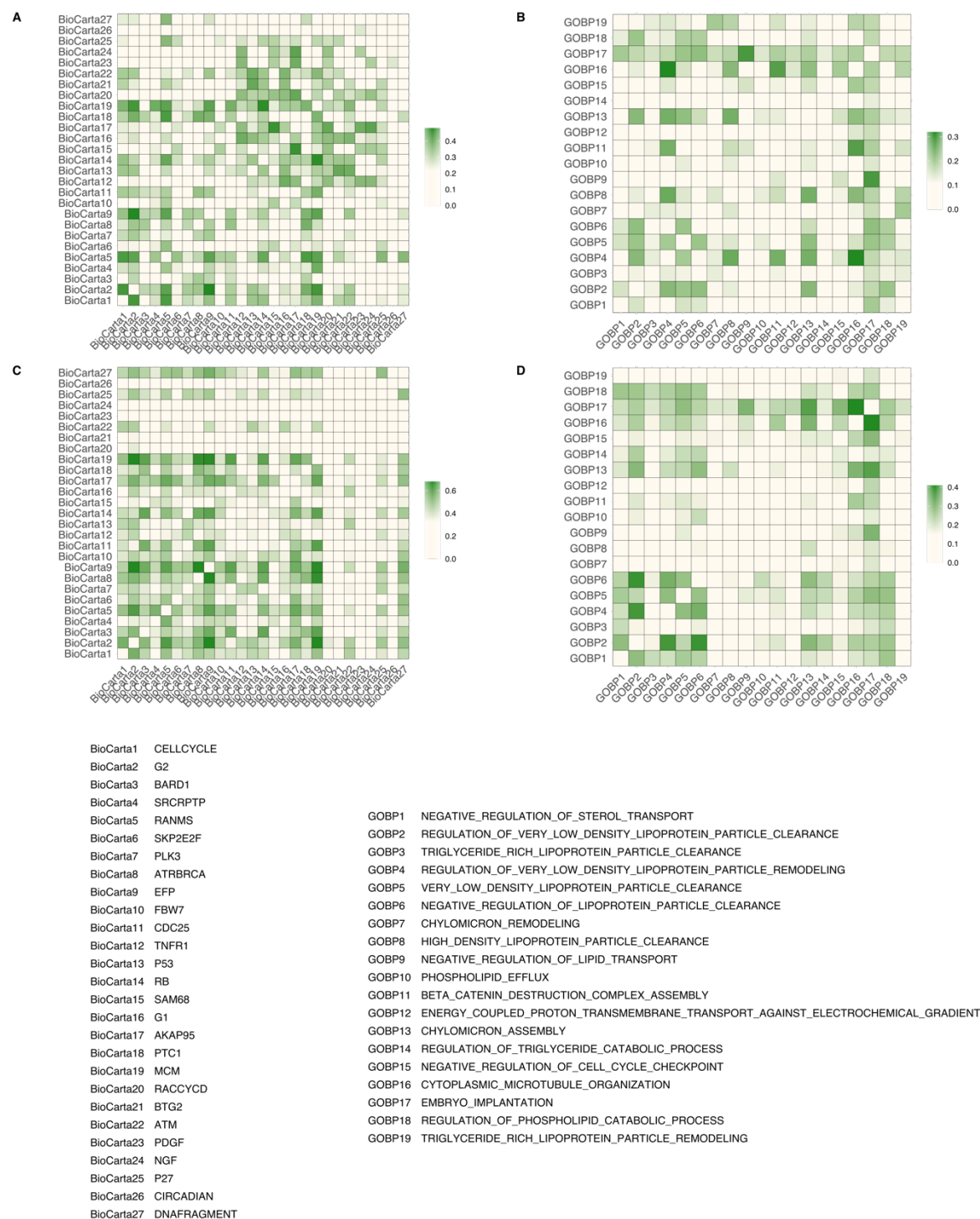


**Figure S18.** Rule-based network displaying the most relevant co-enrichments of annotations obtained from (A, B) the BioCarta collection for the GII vs. GIII model and (C, D) the GOBP collection for the LGG vs. GBM model. The networks show 20 most connected nodes obtained from the top 10% based on the rule connection from a set of significant rules (FDR-adjusted  $P$  value < 0.01). Connection values of nodes and edges represent a strength of co-enrichment from the classifier. Subnetworks were generated separately with respect to the decision class for each RBM.

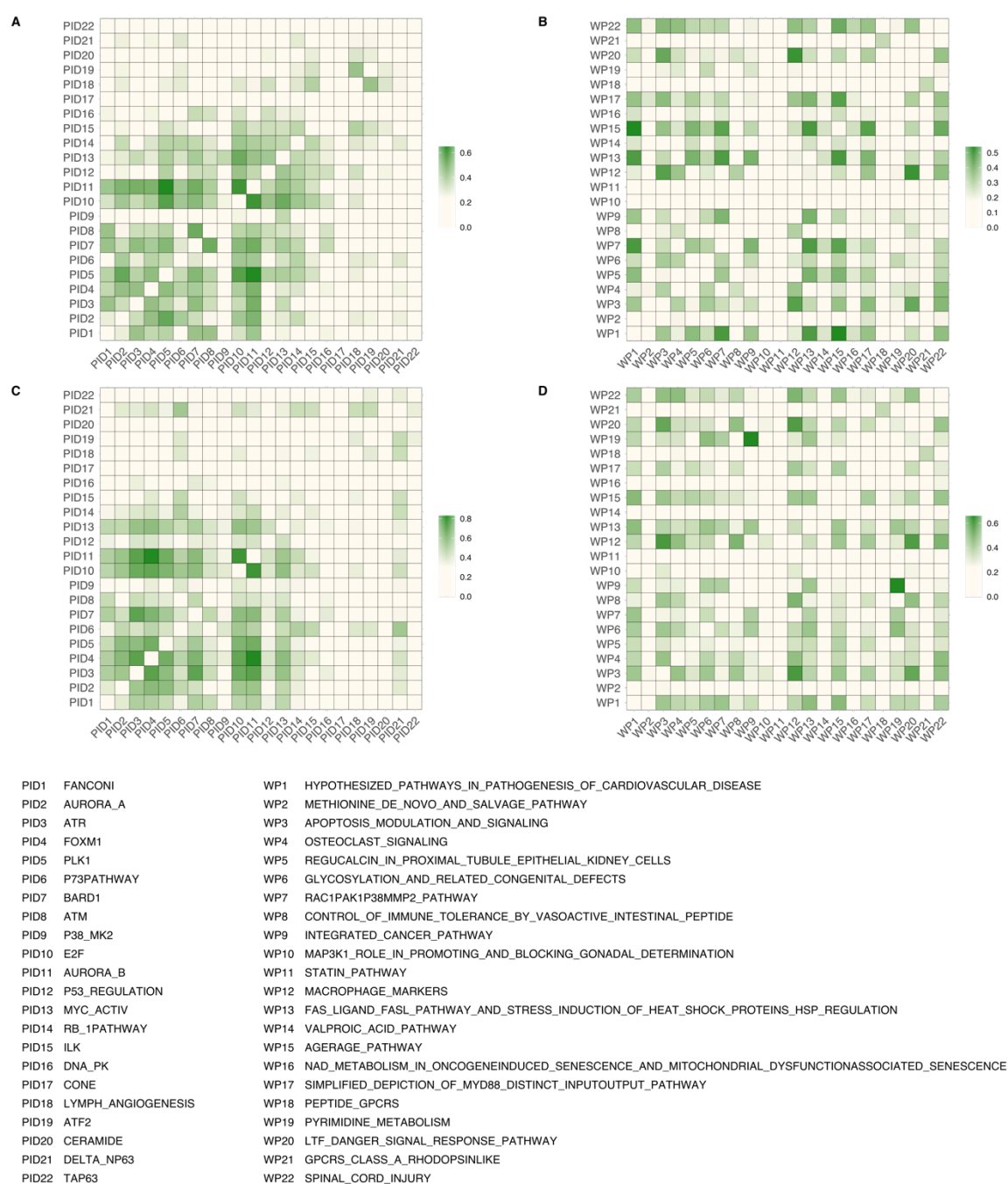




**Figure S20.** Validation of top MSigDB collections: CGP for GII vs. GIII and GOCC for LGG vs. GBM. Heatmaps were generated based on cohort (A, B) CGGA batch 1 and (C, D) CGGA batch 2. Values of total correlation among two variables and a decision class given in nats.



**Figure S21.** Validation of top MSigDB collections: BioCarta pathways for GII vs. GIII and GOBP for LGG vs. GBM. Heatmaps were generated based on cohort (A, B) CGGA batch 1 and (C, D) CGGA batch 2. Values of total correlation among two variables and a decision class given in nats.



**Figure S22.** Validation of top MSigDB collections: PID for GII vs. GIII and WP for LGG vs. GBM. Heatmaps were generated based on cohort (A, B) CGGA batch 1 and (C, D) CGGA batch 2. Values of total correlation among two variables and a decision class given in nats.



## Supplementary Captions for Tables S1-S8

**Table S1.** Computational methods and R packages that were used in this study.

**Table S2.** GM1 and GM2 subsets retrieved from GM modelling of PC1.

**Table S3.** A list of significant DEGs for G2 vs. G3 (FDR < 0.001).

**Table S4.** A list of significant DEGs for LGG vs. GBM (FDR < 0.001).

**Table S5.** Significant enrichment results (FDR < 0.05) for G2 vs. G3 DEGs from *gProfiler*.

**Table S6.** Significant enrichment results (FDR < 0.05) for LGG vs. GBM DEGs from *gProfiler*.

**Table S7.** A list of significant rules (FDR < 0.01) for G2 vs. G3 RBM built for the CGP collection.

**Table S8.** A list of significant rules (FDR < 0.01) for LGG vs. GBM RBM built for the GOCC collection.

## Supplementary References for Table S1

- Draminski, M., & Koronacki, J. (2018). rmcfs: An R Package for Monte Carlo Feature Selection and Interdependency Discovery. *Journal of Statistical Software*, 85, 1–28.  
<https://doi.org/10.18637/JSS.V085.I12>
- Garbulowski, M., Diamanti, K., Smolińska, K., Baltzer, N., Stoll, P., Bornelöv, S., Øhrn, A., Feuk, L., & Komorowski, J. (2021). R.ROSETTA: an interpretable machine learning framework. *BMC Bioinformatics*, 22(1), 1–18. <https://doi.org/10.1186/S12859-021-04049-Z/FIGURES/3>
- Gentleman, R., Carey VJ, Huber W, & Hahne F. (2021). *genefilter: genefilter: methods for filtering* <https://bioconductor.org/packages/release/bioc/html/genefilter.html>
- Hansen, K. D. (2016). IlluminaHumanMethylation450kanno.ilmn12.hg19: Annotation for Illumina's 450k methylation arrays. *R Package Version 0.6.0*, 10, B9.
- Hänzelmann, S., Castelo, R., & Guinney, J. (2013). GSEA: Gene set variation analysis for microarray and RNA-Seq data. *BMC Bioinformatics*, 14(1), 1–15. <https://doi.org/10.1186/1471-2105-14-7/FIGURES/7>
- Hornik, K., Buchta, C., & Zeileis, A. (2008). Open-source machine learning: R meets Weka. *Computational Statistics* 24:2, 24(2), 225–232. <https://doi.org/10.1007/S00180-008-0119-7>
- John, C. R., Watson, D., Russ, D., Goldmann, K., Ehrenstein, M., Pitzalis, C., Lewis, M., & Barnes, M. (2020). M3C: Monte Carlo reference-based consensus clustering. *Scientific Reports* 2020 10:1, 10(1), 1–14. <https://doi.org/10.1038/s41598-020-58766-1>
- Kassambara, A., Kosinski, M., Biecek, P., & Fabian, S. (2021). Package 'survminer'-drawing survival curves using 'ggplot2.' *R Cran*.
- Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E., & Storey, J. D. (2012). The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics*, 28(6), 882. <https://doi.org/10.1093/BIOINFORMATICS/BTS034>
- Meyer, P. E., & Meyer, M. P. E. (2009). Package 'infotheo.' *R Package Version; Citeseer: Princeton, NJ, USA*, 1.
- Scrucca, L., Fop, M., Murphy, T. B., & Raftery, A. E. (2016). mclust 5: Clustering, Classification and Density Estimation Using Gaussian Finite Mixture Models. *The R Journal*, 8(1), 289. <https://doi.org/10.32614/rj-2016-021>
- Smyth, G. K. (2005). Limma: linear models for microarray data. In *Bioinformatics and computational biology solutions using R and Bioconductor* (pp. 397–420). Springer.
- Therneau, T. M., & Lumley, T. (2015). Package 'survival.' *R Top Doc*, 128(10), 28–33.
- Vu, V. Q. (2011). ggbiplot: A ggplot2 based biplot. *R Package Version 0.55*, 755.