**Supplementary Materials for**

# EigenRF: an improved metabolomics normalization method with scores for reproducibility evaluation on importance rankings of different metabolites

Chencheng Tang[1,2], Dongfang Huang[1], Xudong Xing[1,*] and Hua Yang[1,*]

[1]State Key Laboratory of Natural Medicines, China Pharmaceutical University, Nanjing 211198, China
[2]School of Science, China Pharmaceutical University, Nanjing 211198, China

*Corresponding author. State Key Laboratory of Natural Medicines, China Pharmaceutical University, Nanjing 211198, China. E-mail: yanghuacpu@126.com (H.Y.)

**Table S1** The information of ten other normalization methods.

| Type of method | Name | Annotation | Implementation with R package | Description |
|---|---|---|---|---|
| Method based on QC samples | RLSC | Robust locally weighted scatter plot smoothing.[23] | NormalyzeMets[3] | Correct for batch effects by applying robust locally weighted scatter plot smoothing to quality control samples. |
| | RSC | Robust spline correction.[24] | pmp | Fit a smoothing spline to quality control samples and apply it to normalize the data across different batches. |
| | SVR | Support vector regression. | MetNormalizer[25] | Apply support vector regression to model and correct for intensity drift in metabolomics data based on quality control samples. |
| | SERRF | Systematic error elimination employing random forest. | SERRF[17] | Utilize random forest to model and correct for systematic errors in large-scale untargeted lipidomics data based on quality control samples. |
| Method based on Internal Standard | SIS | Single internal standard.[26] | | Apply a single internal standard to normalize data, assuming the standard's response is indicative of overall technical variability. |
| | NOMIS | Optimal normalization factor based on multiple internal standards.[27] | CCMN[28] | Determine the optimal normalization factor for each metabolite by utilizing the variability information from multiple internal standards and their correlation to the metabolites. |
| | CCMN | Cross-contribution compensation based on multiple internal standards.[28] | | Use multiple internal standards, and apply linear regression and PCA to account for cross-contribution effects. |
| Method based on biological samples | ber | Linear fitting based on position/scale. | ber[6] | Employ a two-stage regression approach to adjust for location and scale batch effects. |
| | WaveICA | Wavelet transform algorithm based on independent component analysis. | WaveICA[19] | Utilize wavelet transform and independent component analysis to remove batch effects by decomposing data into multi-scale components and eliminating non-biological variations. |

| Combination method | ISWSVR | Based on the combination of B-MIS and SVR. | Norm-ISWSVR[4] | Integrate internal standard normalization with support vector regression to remove systematic errors. |
| --- | --- | --- | --- | --- |

**Note S1** The information of the other evaluation metrics.

a)    RSD of QCs

The relative standard deviation (RSD) of features in QC samples is a commonly used metric to evaluate the reproducibility of biomarker discovery studies. Here, we use RSD to represent RSD of QCs. In theory, the signal value of a feature on all QC samples should be the same. The higher the RSD, the more unwanted variation and the lower the reproducibility. It is generally accepted that the number of features with RSD less than 0.3 should account for at least 70% of the total number of features. Additionally, the median RSD of each batch is required to account for systematic errors due to batch effects. The ability of the normalization method to eliminate systematic errors can be evaluated by comparing the median RSD of each batch before and after normalization. The RSD of the $i$-th feature is calculated as follows:

$$RSD_i = \frac{S_i^{qc}}{\overline{y}_i^{qc}},$$

(13)

where $RSD_i$, $S_i^{qc}$ and $\overline{y}_i^{qc}$ denote the relative standard deviation, standard deviation, and mean of the $i$-th feature in the QC samples, respectively.

b)    Run plots

The run plot of the feature is represented by a scatter plot, with the signal values as the vertical axis and the injection sequence or injection time as the horizontal axis. The run plot of the internal standards and non-differential metabolites can effectively visualize the presence of batch effects and signal drift. In theory, the signal values of an internal standard should remain consistent across all samples. However, the batch effects may cause variations between batches, while signal drift may result in a discernible trend along the injection sequence.

c)    PCA plots

Principal component analysis (PCA) is performed on all features to generate PCA plots.

Batch-based PCA plots: The sample points are color-coded based on batch information, illustrating the batch effect through the clustering of samples from the same batch.

Group-based PCA plots: The sample points are color-coded based on classification information, allowing for the assessment of normalization methods by observing alterations in the clustering of samples within each category before and after normalization.

d)    OPLS-DA score plots

Orthogonal partial least squares discriminant analysis is performed on all features of biological samples corresponding to specific categories, resulting in the generation of score plots. The performance of the normalization method can be evaluated by examining the classification shifts of various sample categories before and after normalization.

e)    ROC and PR curves

The data are stratified and then randomly sampled at a ratio of 7:3 to be divided into a training set and a test set. For the training set, differential metabolites are screened according to the criteria of VIP $\geq$ 1, which is derived from the OPLS-DA model, and the FDR-adjusted $p$-value $<$ 0.05, which is obtained from the Mann-Whitney U test. A classification model is subsequently constructed with differential metabolites as predictors and grouping information as response.
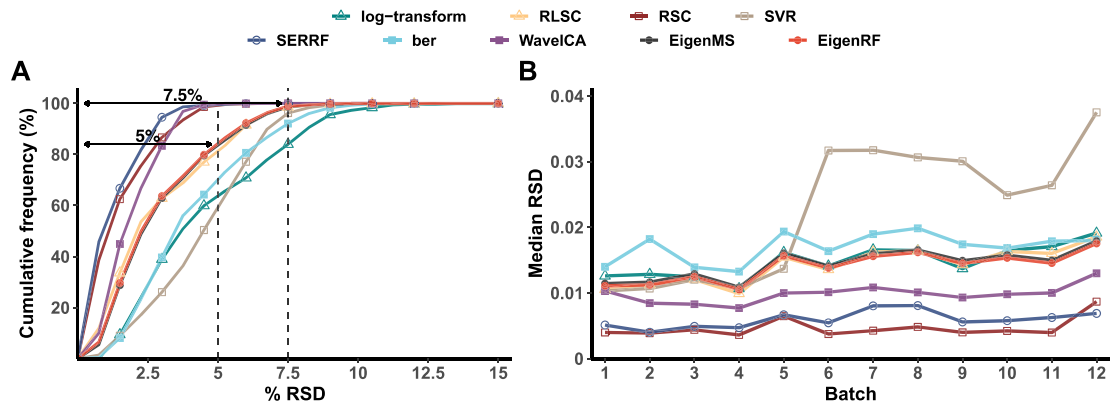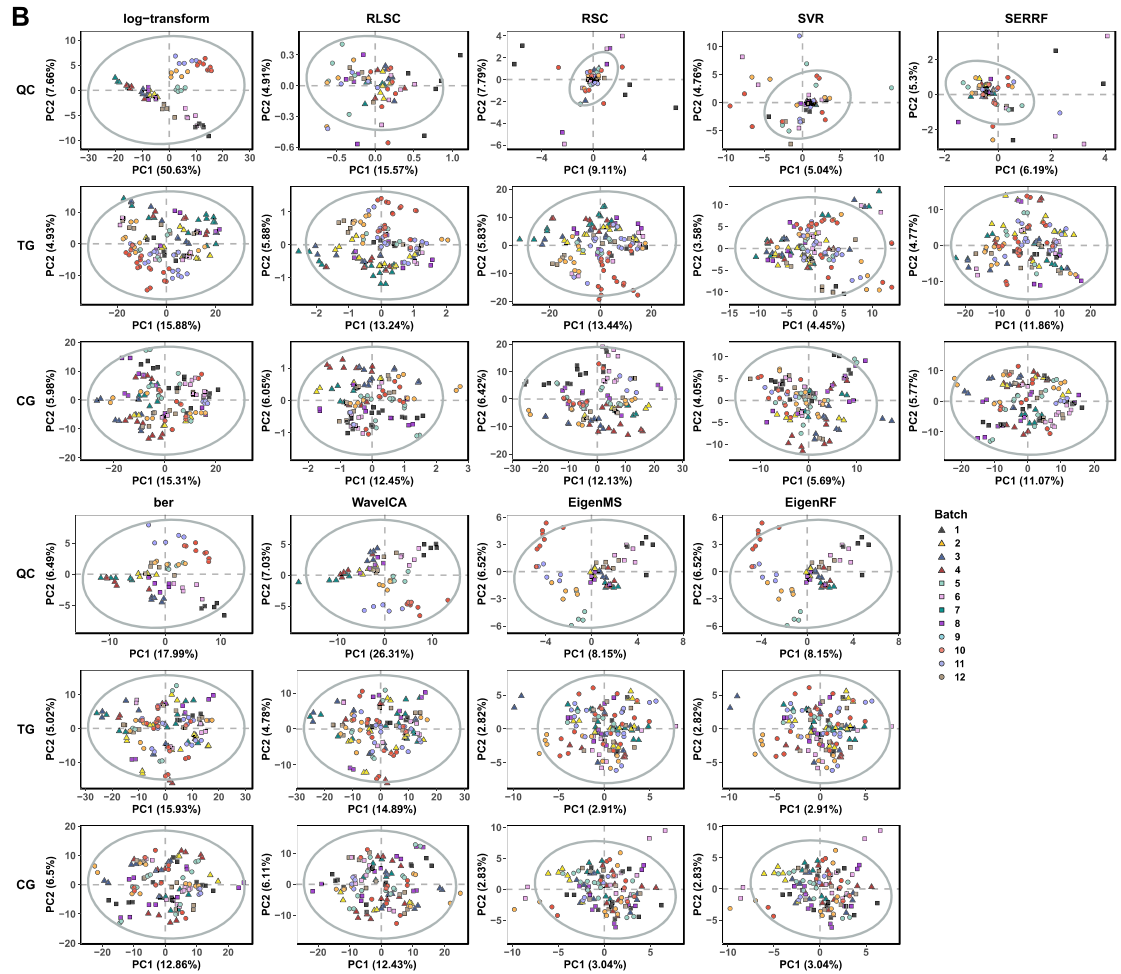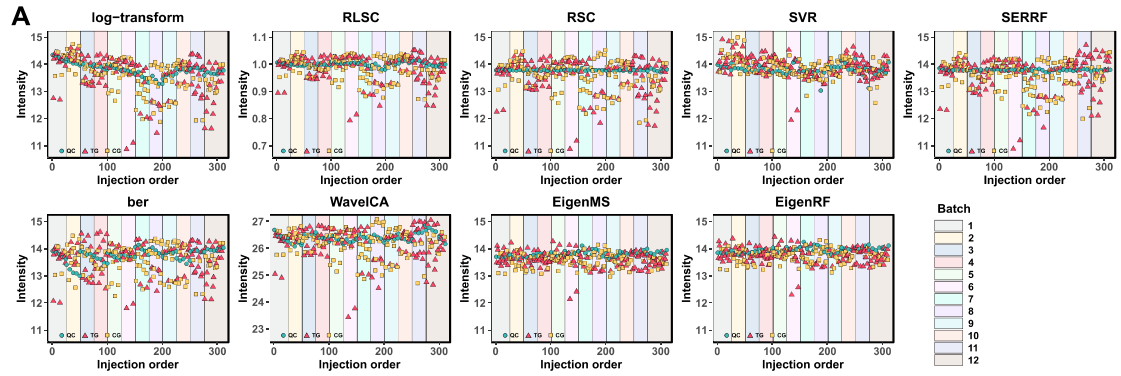
**Fig. S1** Cumulative frequency distribution of the RSDs of QC samples (A) and the median RSD of each batch of QC samples (B) in the BCPUM dataset.
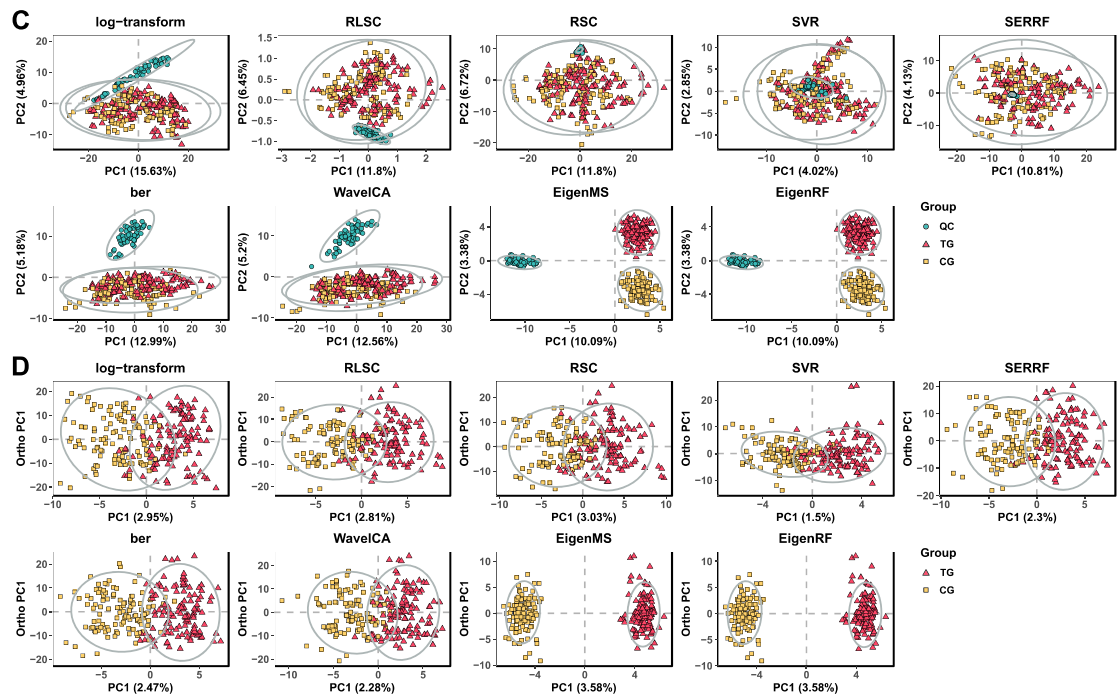
**Fig. S2** Run plots of the first feature (A), batch-based PCA plots of the QC, TG, and CG samples (B), group-based PCA plots (C), and group-based OPLS-DA score plots (D) in the BCPUM dataset.
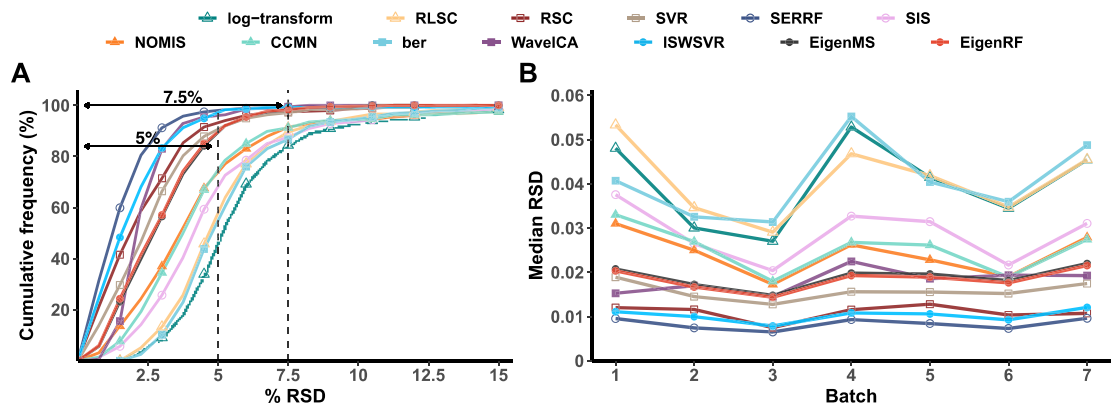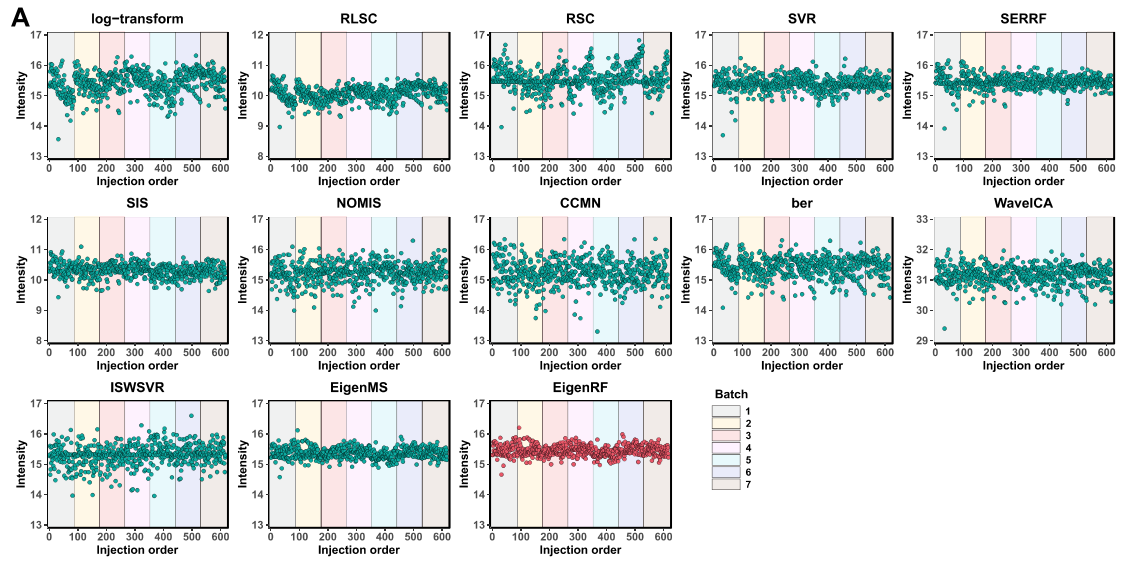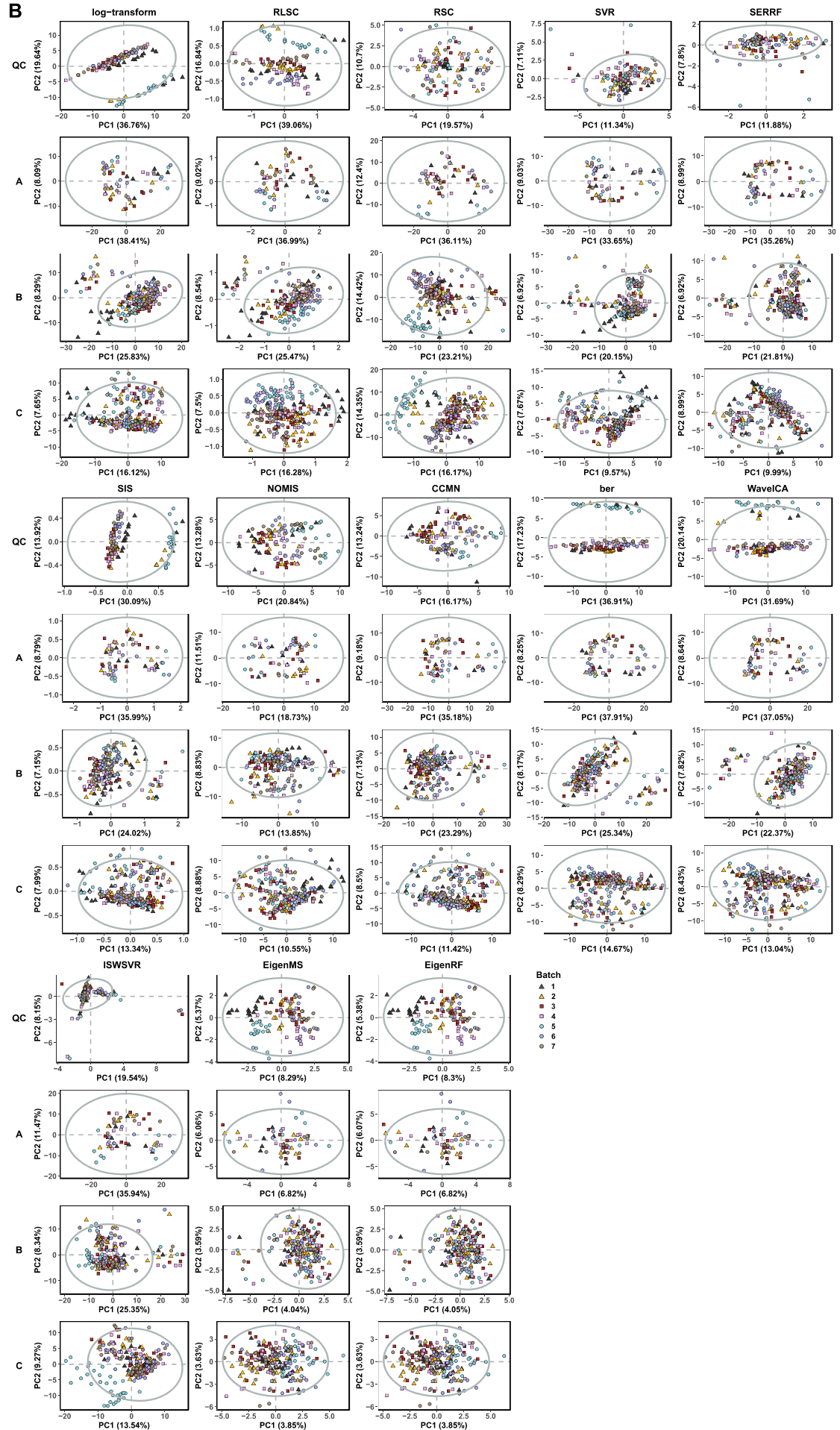
**Fig. S3** Cumulative frequency distribution of the RSDs of QC samples (A) and the median RSD of each batch of QC samples (B) in the GCPPM dataset.

**A**

log-transform | RLSC | RSC | SVR | SERRF

SIS | NOMIS | CCMN | ber | WaveICA

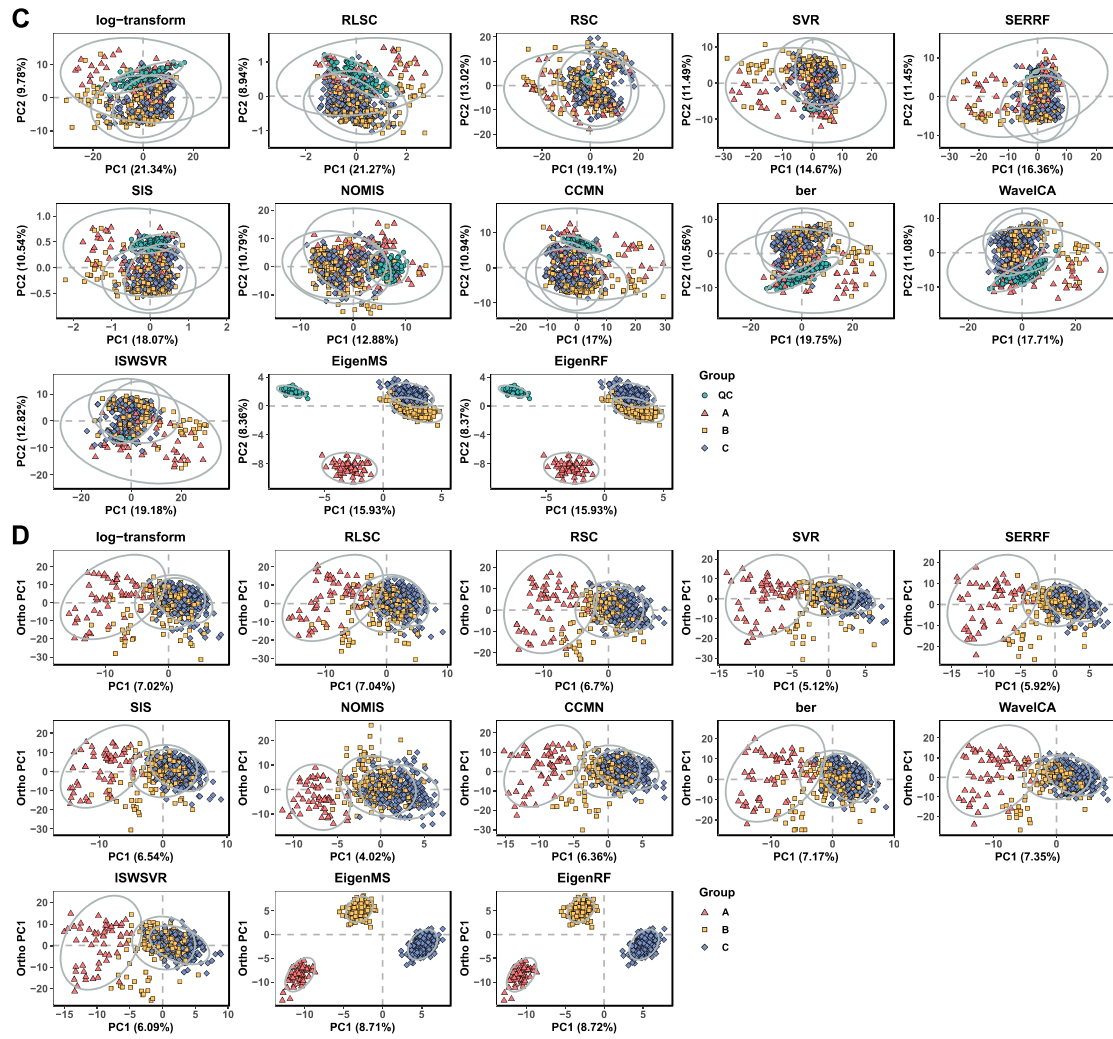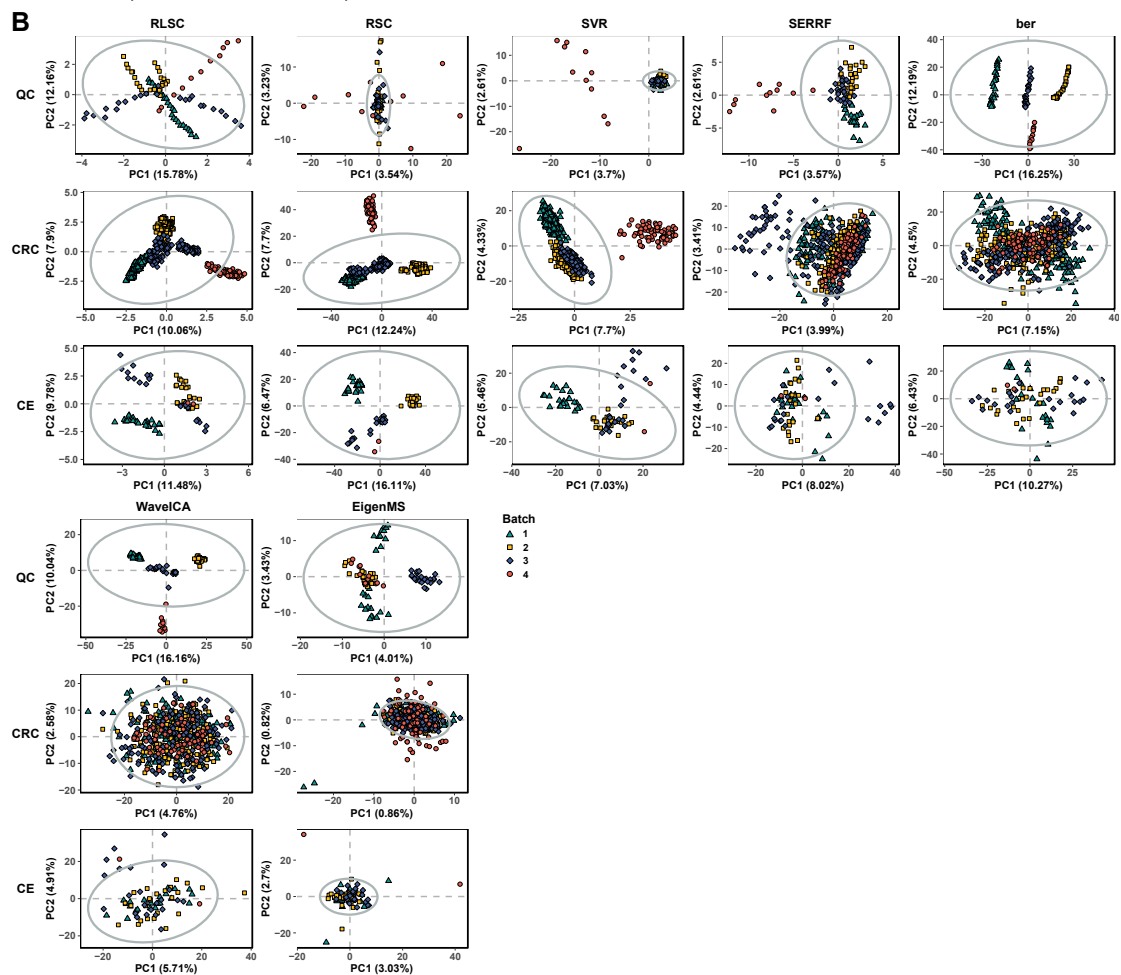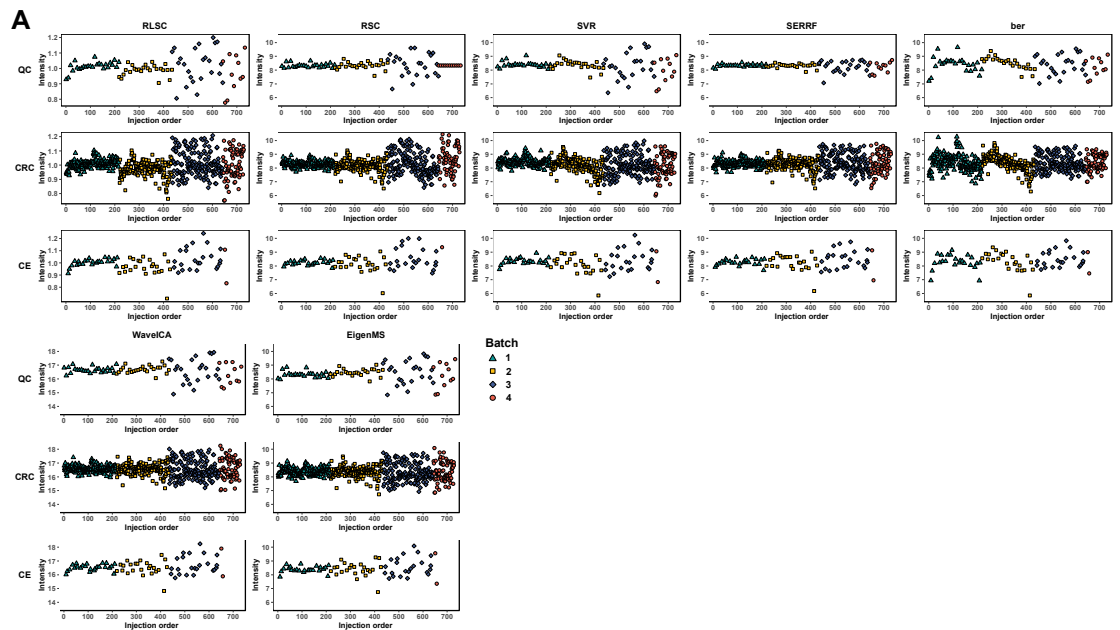ISWSVR | EigenMS | EigenRF

Batch
1
2
3
4
5
6
7

**Fig. S4** Run plots of the internal standard for all samples (A), batch-based PCA plots of the QC, A, B, and C samples (B), group-based PCA plots (C), and group-based OPLS-DA score plots (D) in the GCPPM dataset.
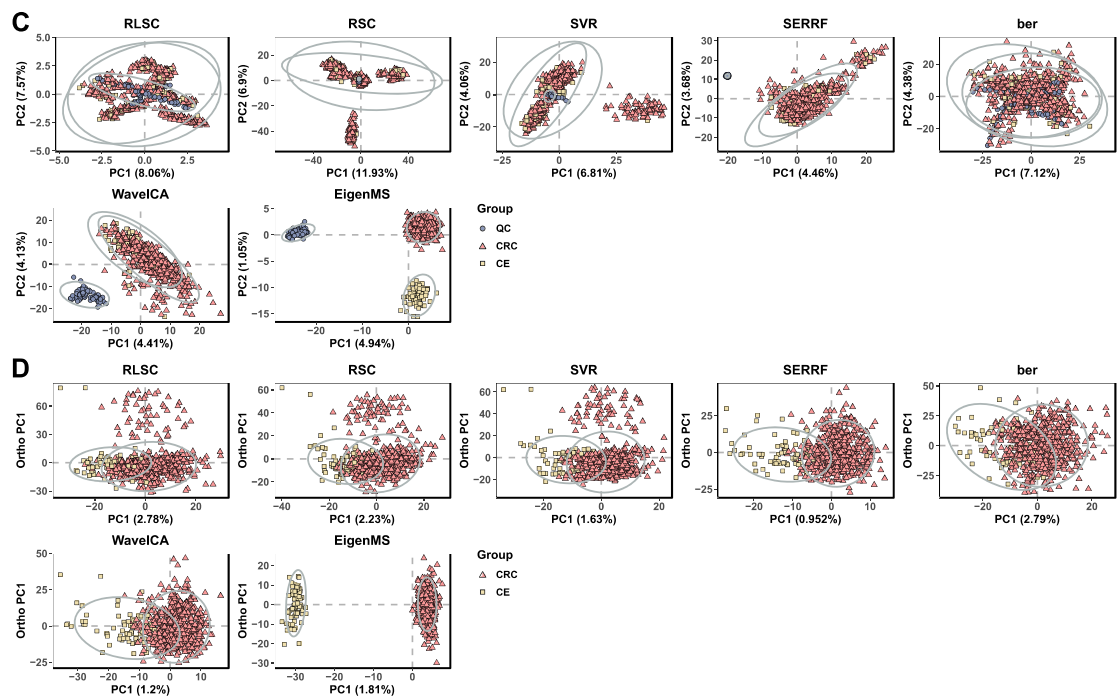
**Fig. S5** Run plots of the first feature (A), batch-based PCA plots of the QC, CRC, and CE samples (B), group-based PCA plots (C), and group-based OPLS-DA score plots (D) in the ACPPM dataset.
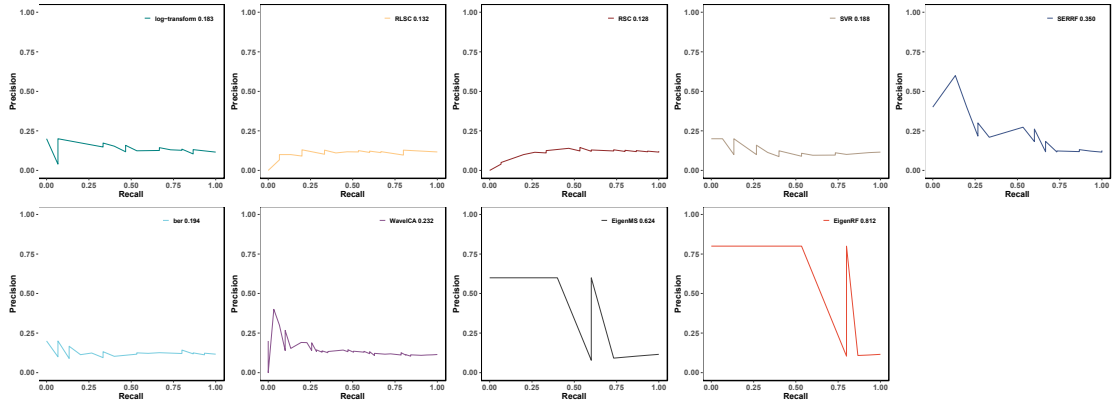
**Fig. S6** PR plots for CE samples of the SVM classification model based on the differential metabolites between the CRC and CE samples in the ACPPM dataset. In the PR plots, the horizontal axis and the vertical axis are respectively the mean recall rate and the mean precision rate of ten-fold cross-validation, and the mean AUC values in the legends.

**References**

23   W.B. Dunn, D. Broadhurst, P. Begley, E. Zelena, S. Francis-McIntyre, N. Anderson, M. Brown, A. Knowles, J.D. Halsall, J.N. Haselden, I.D. Wilson, A.W. Nicholls, D.B. Kell, R. Goodacre and HUSERMET Consortium, *Nat. Protoc.*, 2011, **6**, 1060-1083.

24   J.A. Kirwan, D.I. Broadhurst, R.L. Davidson and M.R. Viant, *Anal. Bioanal. Chem.*, 2013, **405**, 5147-5157.

25   X. Shen, X. Gong, Y. Cai, Y. Guo, J. Tu, H. Li, T. Zhang, J. Wang, F. Xue and Z. Zhu, *Metabolomics*, 2016, **12**, 89.

26   J. Gullberg, P. Jonsson, A. Nordström, M. Sjöström and T. Moritz, *Anal. Biochem.*, 2004, **331**, 283-295.

27   M. Sysi-Aho, M. Katajamaa, L. Yetukuri and M. Orešič, *BMC Bioinformatics*, 2007, **8**, 93.

28   H. Redestig, A. Fukushima, H. Stenlund, T. Moritz, M. Arita, K. Saito and M. Kusano, *Anal. Chem.*, 2009, **81**, 7974-7980.