

Supplementary Material

Zero-inflated quantile rank-score based test (ZIQRank) with application to scRNA-seq differential gene expression analysis

Wodan Ling, Wenfei Zhang, Bin Cheng, Ying Wei

The supplementary material includes proofs of Theorem 1, additional discussions, additional figures, and URLs of data sets used in simulation studies and real data applications.

A Technical proofs

A.1 Proof of Theorem 1(a)

Proof. Define

$$\mathbf{m}(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) = \begin{pmatrix} \tilde{\mathbf{X}}\{\tau - I(Y - \tilde{\mathbf{X}}^\top \boldsymbol{\theta}^Q < 0)\} \\ \mathbf{X}\{I(Y > 0) - \pi(\boldsymbol{\theta}^L, \mathbf{X})\} \end{pmatrix},$$

$$\mathbf{M}(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) = \int \mathbf{m}(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) dP, \quad \mathbf{M}_n(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) = \int \mathbf{m}(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) d\mathbb{P}_n,$$

where $\tilde{\mathbf{X}} = (\tilde{\mathbf{Z}}^\top, \tilde{\mathbf{C}}^\top)^\top$.

By the central limit theorem,

$$\sqrt{n}\mathbf{M}_n(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) \xrightarrow{d} N \left\{ \mathbf{0}, \begin{pmatrix} \tau(1-\tau)\mathbf{D}_0 & \mathbf{0}_{p \times p} \\ \mathbf{0}_{p \times p} & \mathbf{D}_{1,\gamma} \end{pmatrix} \right\},$$

where

$$\mathbf{D}_0 = \mathbb{E} \left\{ \tilde{\mathbf{X}}_i \tilde{\mathbf{X}}_i^\top \right\} = \begin{pmatrix} \mathbb{E} \left\{ \tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^\top \right\} & \mathbb{E} \left\{ \tilde{\mathbf{Z}}_i \tilde{\mathbf{C}}_i^\top \right\} \\ \mathbb{E} \left\{ \tilde{\mathbf{C}}_i \tilde{\mathbf{Z}}_i^\top \right\} & \mathbb{E} \left\{ \tilde{\mathbf{C}}_i \tilde{\mathbf{C}}_i^\top \right\} \end{pmatrix} = \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{pmatrix}. \quad (1)$$

Noticing the orthogonality of $\boldsymbol{\theta}^Q(\tau)$ and $\boldsymbol{\theta}^L$ in the asymptotic distribution of $\sqrt{n}\mathbf{M}_n(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L)$ and for simplicity, denote

$$\mathbf{m}_\alpha(\boldsymbol{\theta}^Q(\tau)) = \tilde{\mathbf{Z}}\{\tau - I(Y - \tilde{\mathbf{X}}^\top \boldsymbol{\theta}^Q < 0)\}, \quad \mathbf{m}_\beta(\boldsymbol{\theta}^Q(\tau)) = \tilde{\mathbf{C}}\{\tau - I(Y - \tilde{\mathbf{X}}^\top \boldsymbol{\theta}^Q < 0)\},$$

and

$$\mathbf{M}_{n,\alpha}(\boldsymbol{\theta}^Q(\tau)) = \int \mathbf{m}_\alpha(\boldsymbol{\theta}^Q(\tau)) d\mathbb{P}_n, \quad \mathbf{M}_{n,\beta}(\boldsymbol{\theta}^Q(\tau)) = \int \mathbf{m}_\beta(\boldsymbol{\theta}^Q(\tau)) d\mathbb{P}_n.$$

Since $\hat{\boldsymbol{\theta}}_n^Q(\tau) = (\hat{\boldsymbol{\alpha}}_n(\tau)^\top, \hat{\boldsymbol{\beta}}_n(\tau)^\top)^\top$ is the restricted M-estimator, by Lagrange multi-

plier method,

$$\sqrt{n}\mathbf{M}_{n,\alpha}(\widehat{\boldsymbol{\theta}}_n^Q(\tau)) = \mathbf{0}, \quad (2)$$

$$\sqrt{n}\mathbf{M}_{n,\beta}(\widehat{\boldsymbol{\theta}}_n^Q(\tau)) + \sqrt{n}\widehat{\boldsymbol{\lambda}}_n = \mathbf{0}, \quad (3)$$

where $\widehat{\boldsymbol{\lambda}}_n$ is the estimator of the Lagrange multiplier parameter.

Let $\boldsymbol{\Delta}$ be a compact neighborhood of the true $\boldsymbol{\theta}^L$, T be a compact subset of interval $(0, 1)$ containing $\Gamma(\tau; \mathbf{X}, \boldsymbol{\theta}^L)$, and $\boldsymbol{\Psi}$ is a compact neighborhood of the true $\boldsymbol{\theta}^Q \circ \Gamma(\tau; \mathbf{X}, \boldsymbol{\theta}^L)$, where

$$\boldsymbol{\theta}^Q \circ \Gamma(\tau; \mathbf{X}, \boldsymbol{\theta}^L) = \boldsymbol{\theta}^Q(\tau_s), \text{ and } \tau_s = \Gamma(\tau; \mathbf{X}, \boldsymbol{\theta}^L) = \max\left(\frac{\tau - \{1 - \pi(\boldsymbol{\theta}^L, \mathbf{X})\}}{\pi(\boldsymbol{\theta}^L, \mathbf{X})}, 0\right).$$

The class

$$\mathcal{F} = \left\{ \mathbf{m}(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L), \tau \times \boldsymbol{\theta}^L \times \boldsymbol{\theta}^Q \in T \times \boldsymbol{\Delta} \times \boldsymbol{\Psi} \right\}$$

is clearly a VC class with a squared integrable envelope function $\mathbf{X}\mathbf{X}^\top$. Thus, \mathcal{F} is a Donsker class. Define

$$\mathbb{G}_n(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) = \sqrt{n} \{ \mathbf{M}_n(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) - \mathbf{M}(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) \},$$

The fact that \mathcal{F} is a Donsker class and both $\widehat{\boldsymbol{\theta}}_n^Q(\tau)$ and $\widehat{\boldsymbol{\theta}}_n^L$ are consistent implies that

$$\mathbb{G}_n(\widehat{\boldsymbol{\theta}}_n^Q(\tau), \widehat{\boldsymbol{\theta}}_n^L) = \mathbb{G}_n(\boldsymbol{\theta}^Q(\tau), \boldsymbol{\theta}^L) + o_p(1),$$

which implies

$$\sqrt{n}\mathbf{M}_n(\widehat{\boldsymbol{\theta}}_n^Q(\tau)) = \sqrt{n}\mathbf{M}_n(\boldsymbol{\theta}^Q(\tau)) + \nabla\mathbf{M}(\boldsymbol{\theta}^Q(\tau))\sqrt{n}(\widehat{\boldsymbol{\alpha}}_n(\tau) - \boldsymbol{\alpha}(\tau)) + o_p(1), \quad (4)$$

where,

$$-\nabla\mathbf{M}(\boldsymbol{\theta}^Q(\tau)) = \begin{pmatrix} \mathbb{E} \left\{ f_i(\tilde{\mathbf{X}}_i^\top \boldsymbol{\theta}^Q(\tau)) \tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^\top \right\} & \mathbb{E} \left\{ f_i(\tilde{\mathbf{X}}_i^\top \boldsymbol{\theta}^Q(\tau)) \tilde{\mathbf{Z}}_i \tilde{\mathbf{C}}_i^\top \right\} \\ \mathbb{E} \left\{ f_i(\tilde{\mathbf{X}}_i^\top \boldsymbol{\theta}^Q(\tau)) \tilde{\mathbf{C}}_i \tilde{\mathbf{Z}}_i^\top \right\} & \mathbb{E} \left\{ f_i(\tilde{\mathbf{X}}_i^\top \boldsymbol{\theta}^Q(\tau)) \tilde{\mathbf{C}}_i \tilde{\mathbf{C}}_i^\top \right\} \end{pmatrix} = b \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{pmatrix},$$

with the last equality by Assumption 4.

Plugging (4) into (2) and (3) to get

$$\sqrt{n}\mathbf{M}_{n,\alpha}(\boldsymbol{\theta}^Q(\tau)) - b\mathbf{A}\sqrt{n}(\widehat{\boldsymbol{\alpha}}_n(\tau) - \boldsymbol{\alpha}(\tau)) = o_p(1), \quad (5)$$

$$\sqrt{n}\mathbf{M}_{n,\beta}(\boldsymbol{\theta}^Q(\tau)) - b\mathbf{B}^\top\sqrt{n}(\widehat{\boldsymbol{\alpha}}_n(\tau) - \boldsymbol{\alpha}(\tau)) + \sqrt{n}\widehat{\boldsymbol{\lambda}}_n = o_p(1). \quad (6)$$

Multiply (5) by $-\mathbf{B}^\top \mathbf{A}^{-1}$ from left and add it to (6), then rearrange terms to yield

$$\begin{aligned}
\sqrt{n}\widehat{\boldsymbol{\lambda}}_n &= \sqrt{n} \{ \mathbf{B}^\top \mathbf{A}^{-1} \mathbf{M}_{n,\alpha}(\boldsymbol{\theta}^Q(\tau)) - \mathbf{M}_{n,\beta}(\boldsymbol{\theta}^Q(\tau)) \} + o_p(1) \\
&= (\mathbf{B}^\top \mathbf{A}^{-1}, -\mathbf{I}) \sqrt{n} \begin{pmatrix} \mathbf{M}_{n,\alpha}(\boldsymbol{\theta}^Q(\tau)) \\ \mathbf{M}_{n,\beta}(\boldsymbol{\theta}^Q(\tau)) \end{pmatrix} + o_p(1) \\
&\xrightarrow{d} (\mathbf{B}^\top \mathbf{A}^{-1}, -\mathbf{I}) \cdot N \left\{ \mathbf{0}, \tau(1-\tau) \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{D} \end{pmatrix} \right\} \\
&= N \left(0, \tau(1-\tau)(\mathbf{D} - \mathbf{B}^\top \mathbf{A}^{-1} \mathbf{B}) \right).
\end{aligned}$$

Then, by (3),

$$\check{\mathbf{S}}_{n,\tau}^Q(\widehat{\boldsymbol{\alpha}}_n(\tau)) = \sqrt{n} \mathbf{M}_{n,\beta}(\widehat{\boldsymbol{\theta}}_n^Q(\tau)) = -\sqrt{n}\widehat{\boldsymbol{\lambda}}_n \xrightarrow{d} N \left(0, \tau(1-\tau)(\mathbf{D} - \mathbf{B}^\top \mathbf{A}^{-1} \mathbf{B}) \right),$$

as claimed. □

A.2 Proof of Theorem 1 (b) and (c)

Proof. Replace $\check{\mathbf{C}}_i$ by $\check{\mathbf{C}}_i = \tilde{\mathbf{C}}_i - \mathbb{E}(\tilde{\mathbf{C}}_i | \tilde{\mathbf{Z}}_i)$ to obtain $\check{\mathbf{S}}_{n,\tau}^Q(\widehat{\boldsymbol{\alpha}}_n(\tau))$, and repeat the same argument as Section A.1, then

$$\check{\mathbf{S}}_{n,\tau}^Q(\widehat{\boldsymbol{\alpha}}_n(\tau)) \xrightarrow{d} N(0, \tau(1-\tau)\boldsymbol{\Sigma}_0),$$

where

$$\boldsymbol{\Sigma}_0 = \mathbb{E}(\check{\mathbf{C}}_i \check{\mathbf{C}}_i^\top) - \mathbb{E}(\check{\mathbf{C}}_i \tilde{\mathbf{Z}}_i^\top) \mathbb{E}(\tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^\top)^{-1} \mathbb{E}(\tilde{\mathbf{Z}}_i \check{\mathbf{C}}_i^\top).$$

By definition, $\mathbb{E}(\check{\mathbf{C}}_i \tilde{\mathbf{Z}}_i^\top) = \mathbf{0}$. Hence, $\boldsymbol{\Sigma}_0 = \mathbb{E}(\check{\mathbf{C}}_i \check{\mathbf{C}}_i^\top)$.

Now we need to prove $\mathbf{S}_{n,\tau}^Q(\widehat{\boldsymbol{\alpha}}_n(\tau)) \xrightarrow{d} N(0, \tau(1-\tau)\boldsymbol{\Sigma}_0)$, which is equivalent to prove

$$\mathbf{S}_{n,\tau}^Q(\widehat{\boldsymbol{\alpha}}_n(\tau)) - \check{\mathbf{S}}_{n,\tau}^Q(\widehat{\boldsymbol{\alpha}}_n(\tau)) = o_p(1).$$

Consider the working model

$$\tilde{\mathbf{C}}_i = \boldsymbol{\xi}^\top \tilde{\mathbf{Z}}_i + \boldsymbol{\delta}_i,$$

where $\boldsymbol{\delta}_i$ are i.i.d. random vectors with zero mean and positive definite variance. Then,

$$\tilde{\mathbf{C}}_i^* - \check{\mathbf{C}}_i = -(\widehat{\boldsymbol{\xi}}_n - \boldsymbol{\xi})^\top \tilde{\mathbf{Z}}_i,$$

where $\widehat{\boldsymbol{\xi}}_n$ is the least square estimator of $\boldsymbol{\xi}$. By the consistency of $\widehat{\boldsymbol{\xi}}_n$, which is implied by Assumption 3, and the asymptotic tightness of $\sqrt{n}\mathbf{M}_{n,\alpha}(\boldsymbol{\theta}^Q(\tau))$ established in Section A.1, we have

$$\mathbf{S}_{n,\tau}^Q(\widehat{\boldsymbol{\alpha}}_n(\tau)) - \check{\mathbf{S}}_{n,\tau}^Q(\widehat{\boldsymbol{\alpha}}_n(\tau)) = -(\widehat{\boldsymbol{\xi}}_n - \boldsymbol{\xi})^\top \{ \sqrt{n}\mathbf{M}_{n,\alpha}(\boldsymbol{\theta}^Q(\tau)) \} = o_p(1)O_p(1) = o_p(1),$$

which completes the proof.

The asymptotic properties of $\mathbf{S}_n^Q = (\mathbf{S}_{n,\tau_1}^Q, \dots, \mathbf{S}_{n,\tau_K}^Q)^\top$ can be derived as an extension of the result on a fixed quantile. As shown in [Koenker \(2005\)](#), the between-quantile correlation is

$$\min\{\tau_i, \tau_j\} - \tau_i\tau_j.$$

Details of this proof are omitted.

Next, since $\tilde{\mathbf{C}}_i^* = \check{\mathbf{C}}_i - (\hat{\boldsymbol{\xi}}_n - \boldsymbol{\xi})^\top \tilde{\mathbf{Z}}_i$, we have

$$\begin{aligned} \frac{\tilde{\mathbf{C}}^{*\top} \tilde{\mathbf{C}}^*}{n} &= \frac{\check{\mathbf{C}}^\top \check{\mathbf{C}}}{n} - 2(\hat{\boldsymbol{\xi}}_n - \boldsymbol{\xi})^\top \frac{\tilde{\mathbf{Z}}^\top \check{\mathbf{C}}}{n} + (\hat{\boldsymbol{\xi}}_n - \boldsymbol{\xi})^\top \frac{\tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}}}{n} (\hat{\boldsymbol{\xi}}_n - \boldsymbol{\xi}) \\ &= \mathbb{E}(\check{\mathbf{C}}_i \check{\mathbf{C}}_i^\top) - 2o_p(1)\mathbb{E}(\tilde{\mathbf{Z}}_i \check{\mathbf{C}}_i^\top) + o_p(1)\mathbb{E}(\tilde{\mathbf{Z}}_i \tilde{\mathbf{Z}}_i^\top) \\ &= \mathbb{E}(\check{\mathbf{C}}_i \check{\mathbf{C}}_i^\top) + o_p(1), \end{aligned}$$

which implies that

$$n^{-1} \tilde{\mathbf{C}}^{*\top} \tilde{\mathbf{C}}^* \xrightarrow{p} \mathbb{E}(\check{\mathbf{C}}_i \check{\mathbf{C}}_i^\top) = \boldsymbol{\Sigma}_0.$$

Therefore, by Slutsky's Theorem, given $\boldsymbol{\beta}(\tau) = \mathbf{0}$,

$$T_\tau^Q = \mathbf{S}_{n,\tau}^{Q\top} \mathbf{V}_{n,\tau}^{-1} \mathbf{S}_{n,\tau}^Q \xrightarrow{d} \chi_q^2.$$

Thus, Theorems 1 (b) and (c) are proved. □

B Coefficient-based test

We also tried the coefficient-based test in the quantile regression adjusting for individual zero-inflation. The asymptotic covariance matrix of $\hat{\mathbf{B}}_n = (\hat{\boldsymbol{\beta}}_n(\tau_1), \dots, \hat{\boldsymbol{\beta}}_n(\tau_K))^\top$ can be estimated by bootstrap. In each round, we resample from the entire sample, and estimate the set of quantile coefficients, $\hat{\mathbf{B}}_n$, based on the positive part of the bootstrapped dataset. Such a procedure will introduce the zero-positive uncertainty into the estimation, as expected. With $\hat{\mathbf{B}}_n^{(b)}$, $b = 1, \dots, B$, we can compute the estimated covariance matrix $\hat{\boldsymbol{\Sigma}}_{\text{coef}}$. Note that we use bootstrap to avoid the computational difficulty in the kernel estimation of the conditional local density.

By some simulation experiments, we find that it has an even higher power compared to the ZIQRank test. However it has two drawbacks. (1) The Type I error is sometimes uncontrolled, while the power is just improved marginally. This is because testing based on coefficients is generally unstable, especially at extreme quantiles ([Chen and Wei 2005](#)). (2) Also, it is computationally intensive because of the resampling. Therefore, we still recommend the ZIQRank test for practical use.

C Partial combination procedures

Because of conditional independence, we can first combine the p-values from quantile regression, and further combine it with the p-value from logistic regression by Fisher's combined

probability test. The final test statistic follows a χ_4^2 distribution, and it is robust to the extent of zero-inflation regardless of the combination of the quantile p-values.

References

- Chen, C. and Y. Wei (2005). Computational issues for quantile regression. *Sankhyā: The Indian Journal of Statistics*, 399–417.
- Koenker, R. (2005). Quantile regression. *Cambridge University Press: Cambridge*.

Additional figures

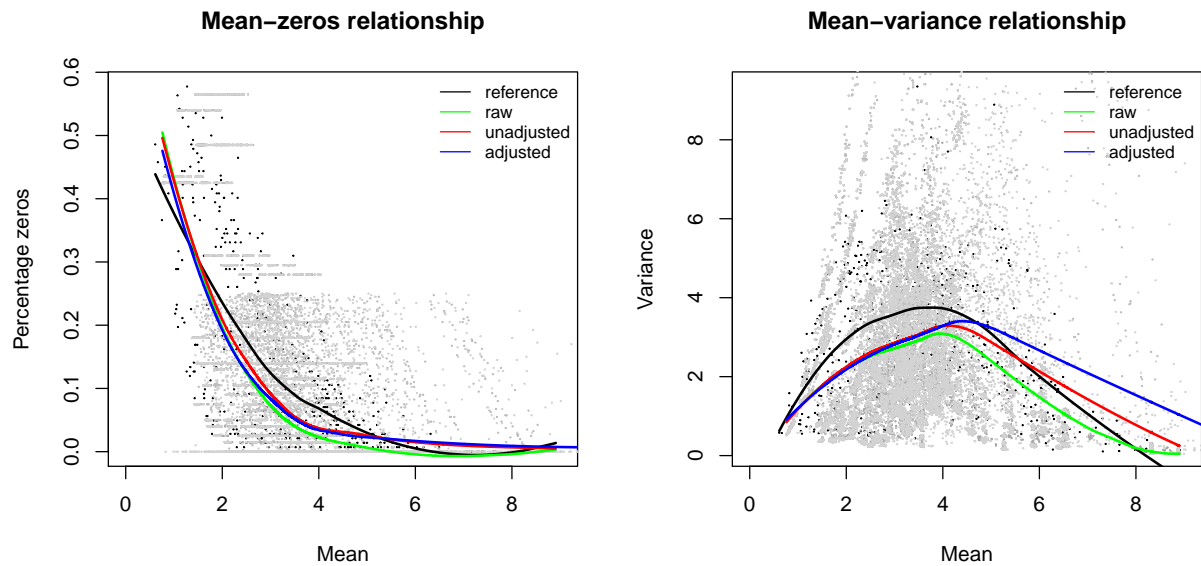
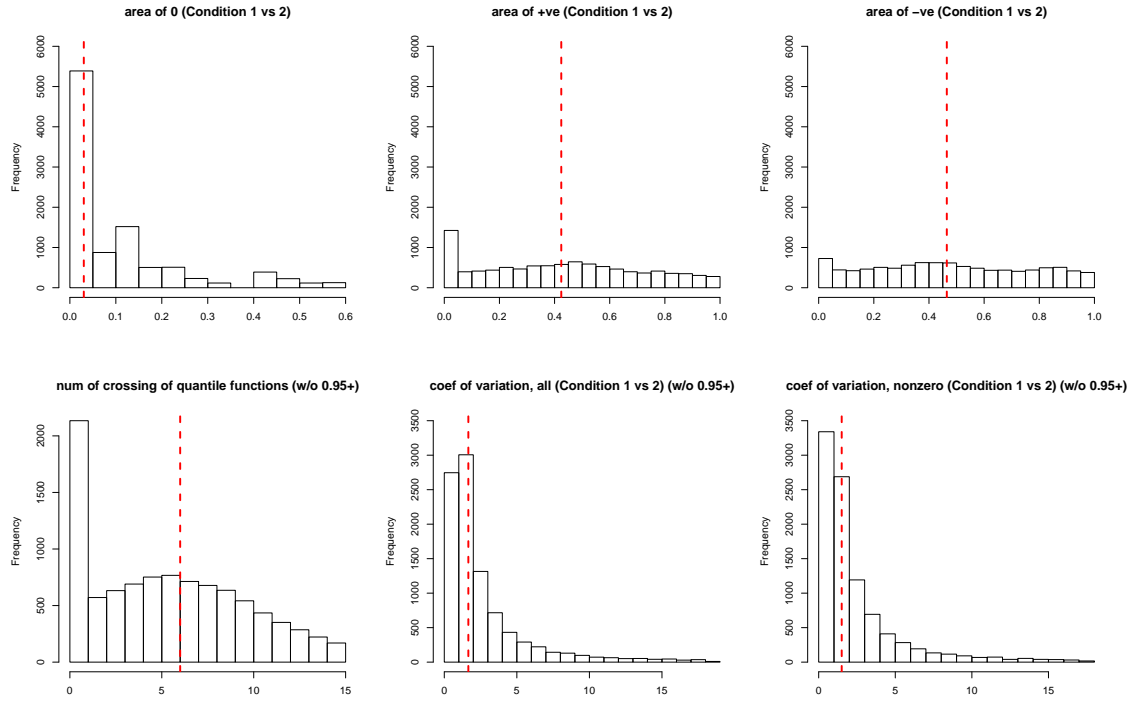
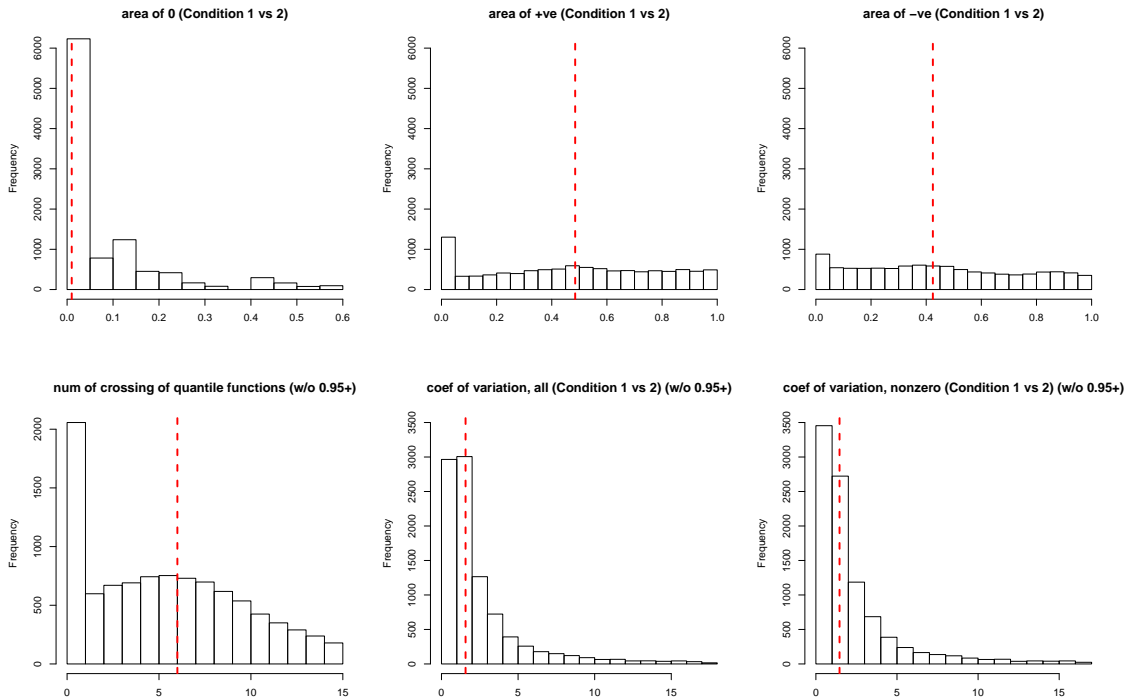


Figure S1: Comparisons of simulated and reference data on the mean-zeros and mean-variance relationships. **reference**: a human embryonic stem cell scRNA-seq data that serves as the starting data of simulation; **raw**: the simulated data by scDD package without modification; **unadjusted**: the simulated data with 0–25% extra zero-inflation; **adjusted**: the simulated data with a continuous confounding covariate and 0–25% extra zero-inflation.

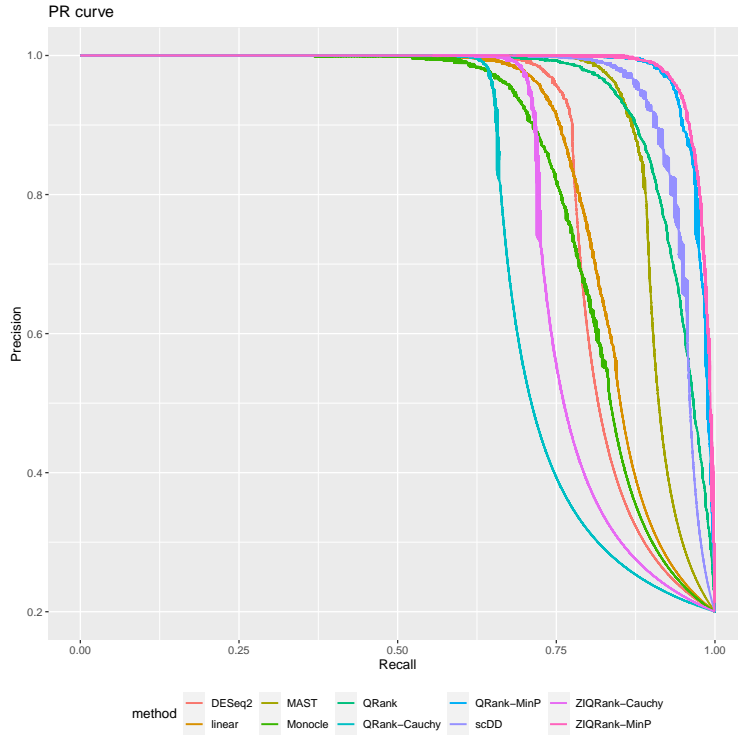


(a) Description of quantile difference in a simulated dataset for unadjusted analysis.

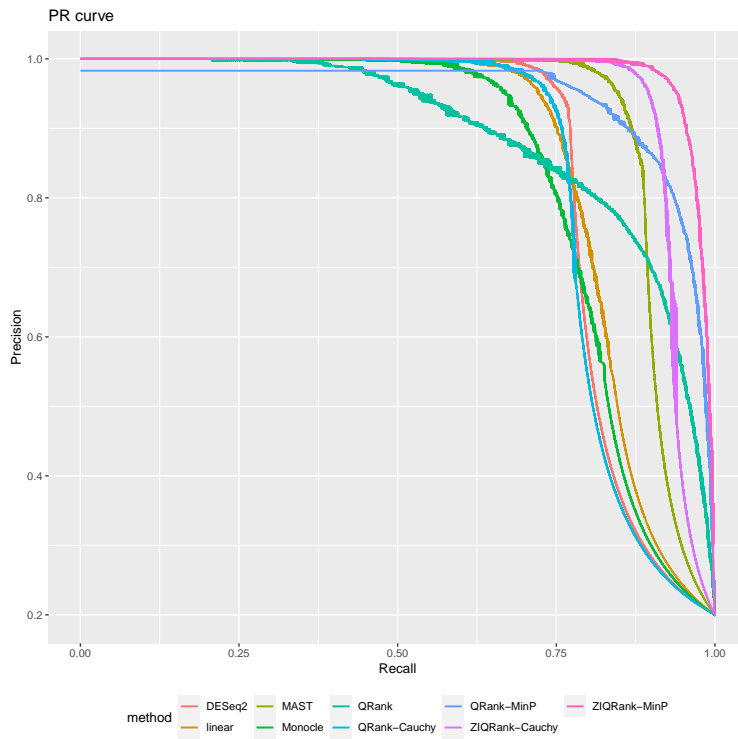


(b) Description of quantile difference in a simulated dataset for adjusted analysis.

Figure S2: Summary of quantile difference between cell conditions in simulated datasets, with median of each statistic marked by dashed red vertical lines.

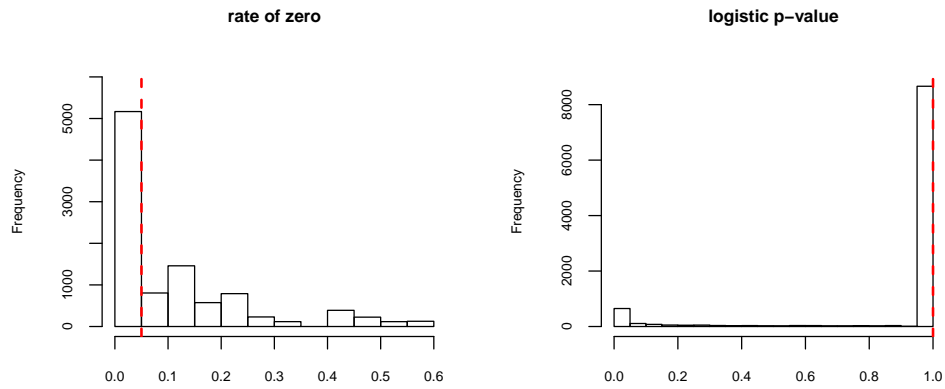


(a) PR curves for all methods in unadjusted simulation study.

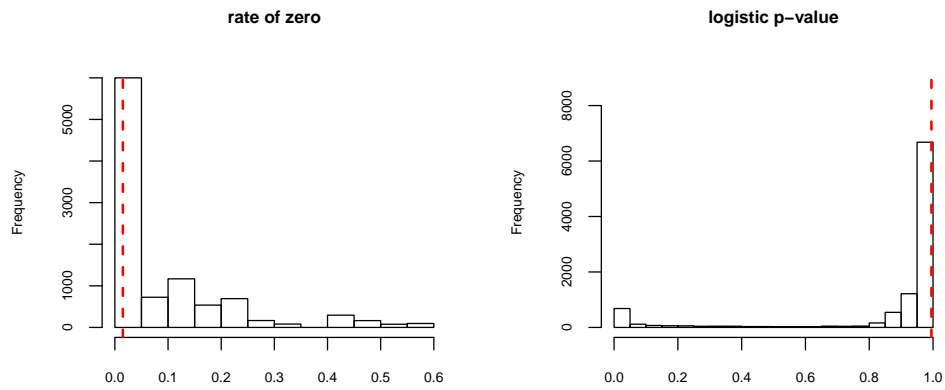


(b) PR curves for all methods in adjusted simulation study.

Figure S3: PR curves of ZIQRank and existing methods in simulation studies.

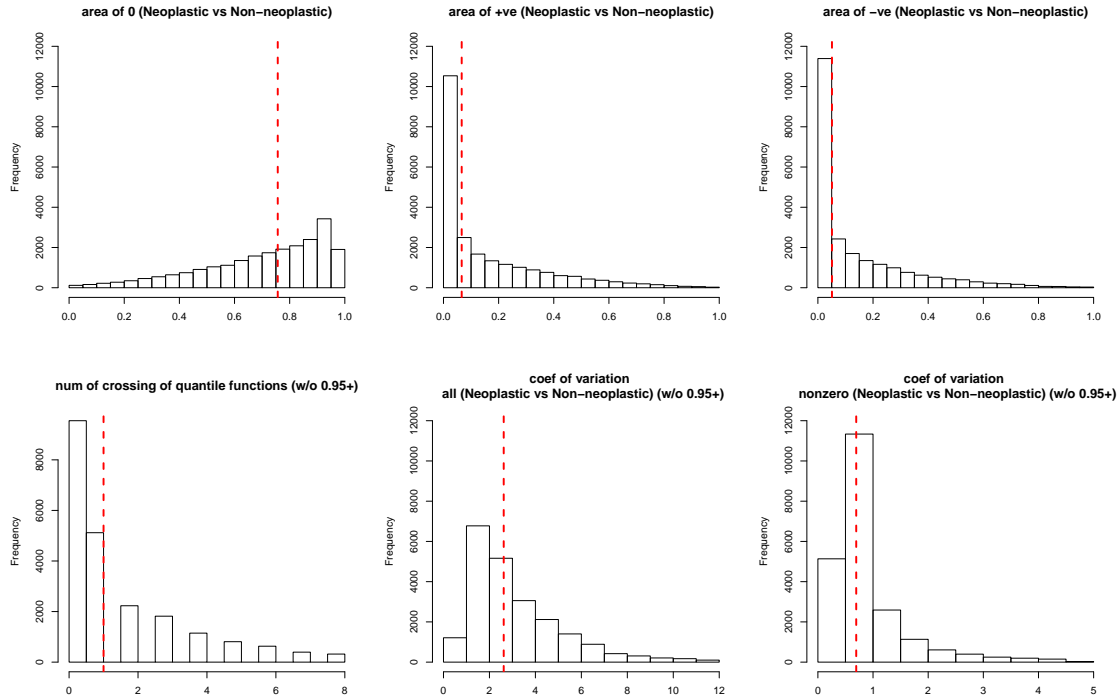


(a) Raw results about zeros on a simulated dataset for unadjusted analysis.

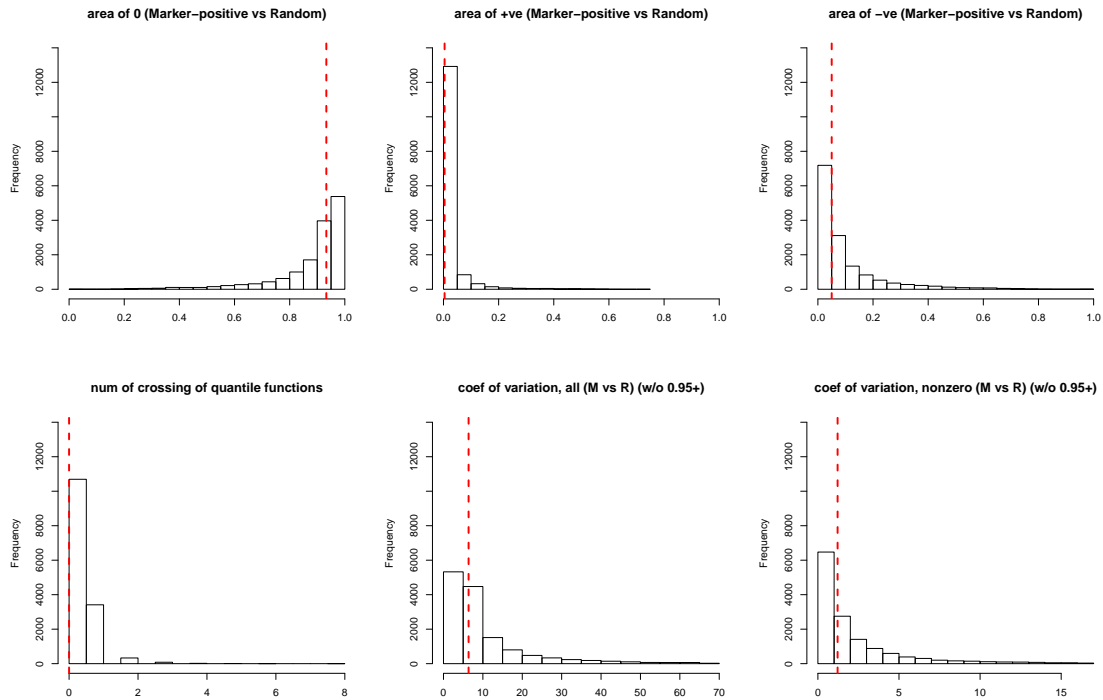


(b) Raw results about zeros on a simulated dataset for adjusted analysis.

Figure S4: Summary of zero rates and p-value from logistic tests in simulated datasets, with median of each statistic marked by dashed red vertical lines.



(a) Summary of quantile difference, neoplastic vs. non-neoplastic cells, in GSE84465.



(b) Summary of quantile difference, marker-positive vs. randomly extracted cells, in GSE62270-GPL17021.

Figure S5: Descriptive statistics about quantile difference between cell conditions in GSE84465 and GSE62270-GPL17021, with median of each statistic marked by dashed red vertical lines.

URLs of data sets

- Simulation starting data, a human embryonic stem cell scRNA-seq data
 - scDatEx, embedded in scDD package: <https://bioconductor.org/packages/release/bioc/html/scDD.html>
- Real data on **conquer** (<http://imlspenticton.uzh.ch:3838/conquer/>)
 - GSE84465 (Homo sapiens, glioblastoma tumors, full-length by Smart-Seq2):
<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE84465>,
<https://pubmed.ncbi.nlm.nih.gov/29091775/>
 - GSE62270-GPL17021 (Mus musculus, cells from mouse intestinal organoids, UMI by CEL-Seq):
<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE62270>,
<https://pubmed.ncbi.nlm.nih.gov/26287467/>