

QRank: Quantile Regression for eQTL Analysis

Gen Li

Department of Biostatistics, Columbia University

March 27, 2017

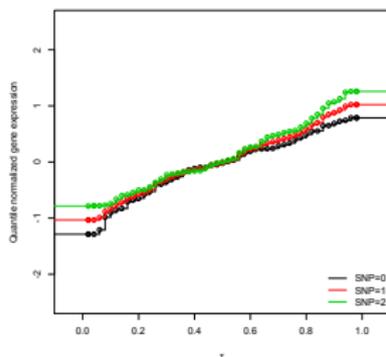
Motivation

- Expression quantitative trait loci (eQTLs) analysis can elucidate genetic regulatory pathways of complex diseases
- Most eQTL studies focus on identifying mean effects
- Linear regression models are commonly used

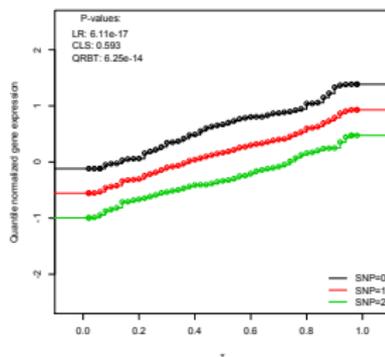
- Expression quantitative trait loci (eQTLs) analysis can elucidate genetic regulatory pathways of complex diseases
- Most eQTL studies focus on identifying mean effects
- Linear regression models are commonly used
- However, genetic variants may affect the entire distribution of gene expressions
- This distributional heterogeneity is understudied

Examples of Distributional Heterogeneity

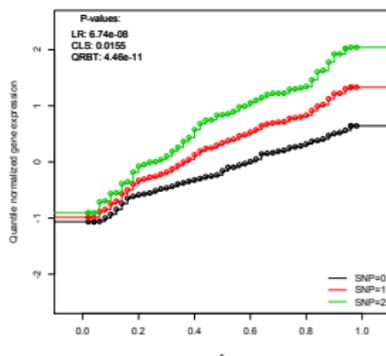
Gene: ENSG00000196932.7 & SNP: 10_63342861_C_T_b37



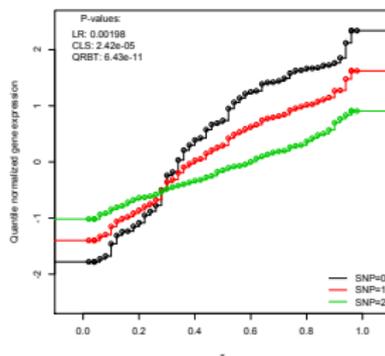
Gene: ENSG00000119943.6 & SNP: 10_100122640_C_G_b37



Gene: ENSG00000259917.1 & SNP: 15_34743556_G_A_b37



Gene: ENSG00000206503.7 & SNP: 6_29880050_T_C_b37



Quantile Regression for eQTL Discovery

In a particular tissue with n samples, for $i \in \{1, 2, \dots, n\}$,

- $Y_{i,k}$ is the expression level of the k th gene in the i th sample
- $x_{i,j}$ is the j th SNP within ± 1 MB of the TSS of the gene in the i th sample
- \mathbf{z}_i is a covariate vector for the i th sample

For each gene-SNP pair k and j , we assume the quantiles of $Y_{i,k}$ follow the model:

$$Q_{Y_{i,k}}(\tau | \mathbf{z}_i, x_{i,j}) = \mathbf{z}_i \alpha_{jk,\tau} + x_{i,j} \beta_{jk,\tau},$$

- **Goal:** identify the k and j pairs for $\beta_{jk,\tau} \neq 0$ for any given $\tau \in (0, 1)$

Define the rank-score function for a fixed quantile τ as

$$S_{n,\tau} = n^{-1/2} \sum_{i=1}^n \rho_{\tau}\{y_{i,k} - \mathbf{z}_i \hat{\alpha}_{jk,\tau}\} x_{i,j}^*$$

- $\rho_{\tau}\{u\} = \tau - \mathbb{I}(u < 0)$ is an asymmetric sign function
- $\hat{\alpha}_{jk,\tau}$ is the estimated coefficient vector under the null $H_0 : \beta_{jk,\tau} = 0$
- $x_{i,j}^*$ is the genotype data adjusted by the covariates
- $S_{n,\tau}$ is close to zero if and only if $\beta_{jk,\tau} = 0$

- Test statistic at a fixed quantile τ :

$$T_{n,\tau} = \mathbf{S}_{n,\tau}^2 / V_n \longrightarrow \chi_1^2, \text{ as } n \rightarrow \infty$$

where V_n^{-1} is the variance of $\mathbf{S}_{n,\tau}$ such that

$$V_n = n^{-1} \tau(1 - \tau) \mathbf{x}_j^{*T} \mathbf{x}_j^*$$

- Composite test statistic $(\tau_1, \dots, \tau_\ell)$:

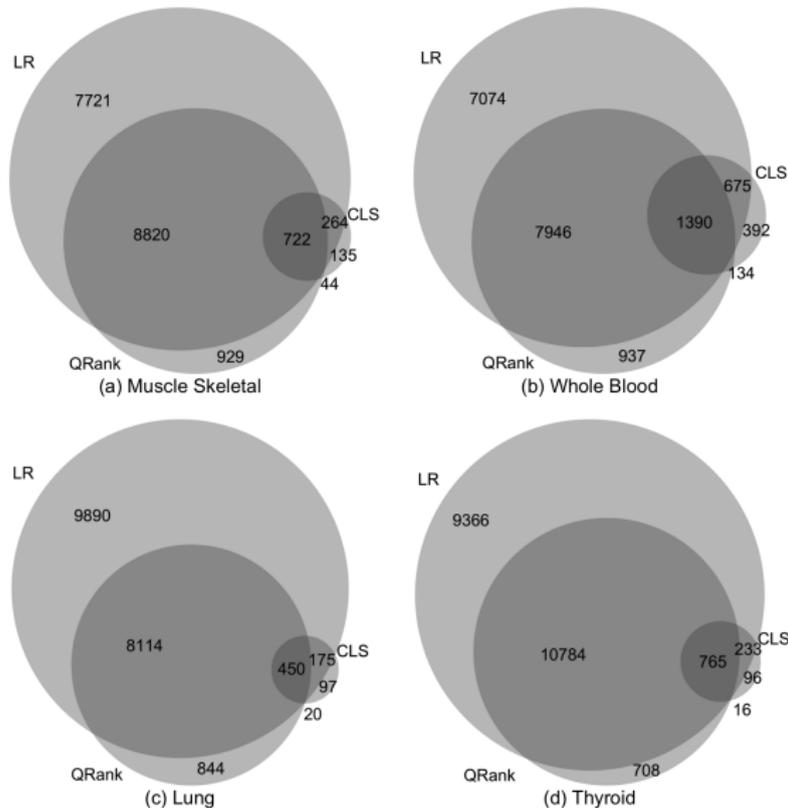
$$T_{n,\ell} = \mathbf{S}_{n,\ell}^T \boldsymbol{\Sigma}_{n,\ell}^{-1} \mathbf{S}_{n,\ell} \longrightarrow \chi_\ell^2, \text{ as } n \rightarrow \infty$$

where $\mathbf{S}_{n,\ell} = (\mathbf{S}_{n,\tau_1}, \dots, \mathbf{S}_{n,\tau_\ell})$ and $\boldsymbol{\Sigma}_{n,\ell}$ has an explicit expression

- Compare QRank (with 5 quantile levels) with LR and CLS¹
- Type I error and power analysis in simulation
- GTEx v6 data analysis (in 4 tissues)
- Gene-level discoveries, tissue specificity, GWAS enrichment

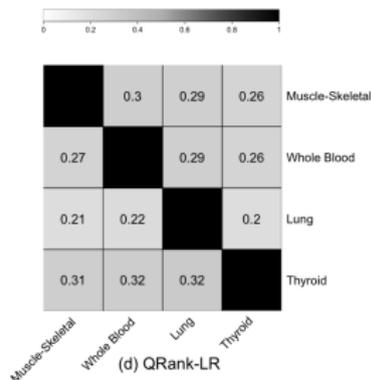
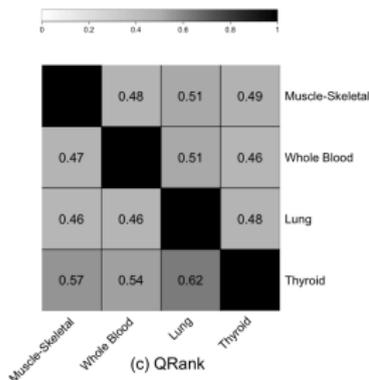
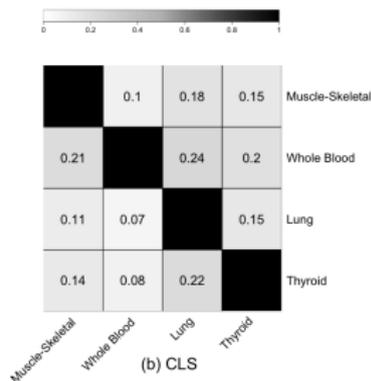
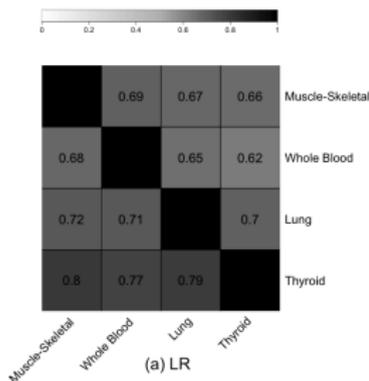
¹Brown et al., *Elife* (2014)

eQTL Discoveries (Unique Genes, FDR=0.05)



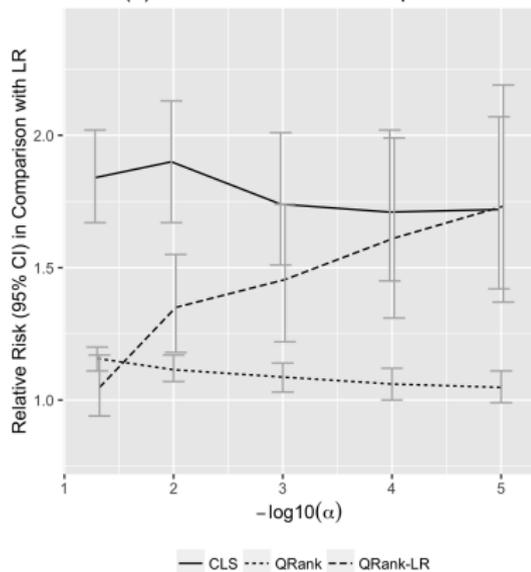
Tissue Specificity

- eQTL sharing coefficient $\pi_{ij} = \mathbb{P}(\text{eQTL in tissue } i \mid \text{eQTL in tissue } j)$

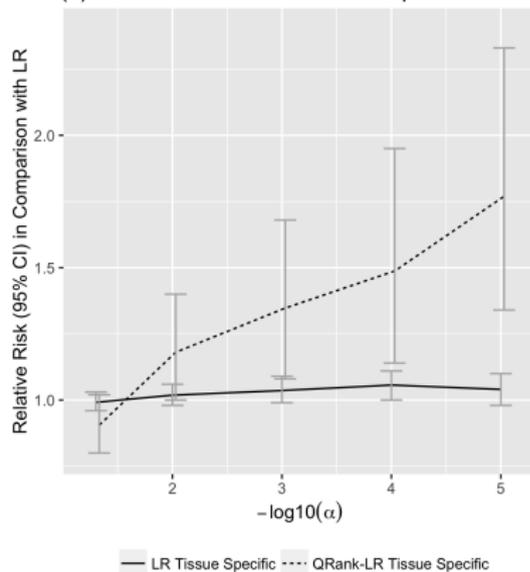


GWAS Enrichment

(a) GWAS Enrichment Comparison



(b) GWAS Enrichment in Tissue-Specific eQTLs



Acknowledgement

Collaborators:

Xiaoyu Song
Zhenwei Zhou
Xianling Wang
Iuliana Ionita-Laza
Ying Wei

Funding:

NSF (DMS-120923)
NIH (R01HG008980)
Calderone Junior Faculty Award
MSPH, Columbia University

Song et al. "QRank: A novel quantile regression tool for eQTL discovery." *Bioinformatics*, 2017+

(<http://biorxiv.org/content/early/2016/08/17/070052>)

R package available at <https://github.com/cran/QRank>