

QCB 508 – Week 12

John D. Storey

Spring 2020

Contents

HD Latent Variable Models	1
Definition	1
Model	1
Estimation	1
Jackstraw	2
Procedure	2
Example: Yeast Cell Cycle	2
Surrogate Variable Analysis	6
Procedure	6
Example: Kidney Expr by Age	7
Extras	10
Source	10
Session Information	11

HD Latent Variable Models

Definition

Latent variables (or hidden variables) are random variables that are present in the underlying probabilistic model of the data, but they are unobserved.

In high-dimensional data, there may be latent variables present that affect many variables simultaneously.

These are latent variables that induce **systematic variation**. A topic of much interest is how to estimate these and incorporate them into further HD inference procedures.

Model

Suppose we have observed data $\mathbf{Y}_{m \times n}$ of m variables with n observations each. Suppose there are r latent variables contained in the r rows of $\mathbf{Z}_{r \times n}$ where

$$\mathbf{E}[\mathbf{Y}_{m \times n} | \mathbf{Z}_{r \times n}] = \mathbf{\Phi}_{m \times r} \mathbf{Z}_{r \times n}.$$

Let's also assume that $m \gg n > r$. The latent variables \mathbf{Z} induce systematic variation in variable \mathbf{y}_i parameterized by ϕ_i for $i = 1, 2, \dots, m$.

Estimation

There exist methods for estimating the row space of \mathbf{Z} with probability 1 as $m \rightarrow \infty$ for a fixed n in two scenarios.

Leek (2011) shows how to do this when $\mathbf{y}_i|\mathbf{Z} \sim \text{MVN}(\phi_i\mathbf{Z}, \sigma_i^2\mathbf{I})$, and the $\mathbf{y}_i|\mathbf{Z}$ are jointly independent.

Chen and Storey (2015) show how to do this when the $\mathbf{y}_i|\mathbf{Z}$ are distributed according to a single parameter exponential family distribution with mean $\phi_i\mathbf{Z}$, and the $\mathbf{y}_i|\mathbf{Z}$ are jointly independent.

Jackstraw

Suppose we have a reasonable method for estimating \mathbf{Z} in the model

$$\mathbb{E}[\mathbf{Y} | \mathbf{Z}] = \Phi \mathbf{Z}.$$

The **jackstraw** method allows us to perform hypothesis tests of the form

$$H_0 : \phi_i = \mathbf{0} \text{ vs } H_1 : \phi_i \neq \mathbf{0}.$$

We can also perform this hypothesis test on any subset of the columns of Φ .

This is a challenging problem because we have to “double dip” in the data \mathbf{Y} , first to estimate \mathbf{Z} , and second to perform significance tests on Φ .

Procedure

The first step is to form estimate $\hat{\mathbf{Z}}$ and then test statistic t_i that performs the hypothesis test for each ϕ_i from \mathbf{y}_i and $\hat{\mathbf{Z}}$ ($i = 1, \dots, m$). Assume that the larger t_i is, the more evidence there is against the null hypothesis in favor of the alternative.

Next we randomly select s rows of \mathbf{Y} and permute them to create data set \mathbf{Y}^0 . Let this set of s variables be indexed by \mathcal{S} . This breaks the relationship between \mathbf{y}_i and \mathbf{Z} , thereby inducing a true H_0 , for each $i \in \mathcal{S}$.

We estimate $\hat{\mathbf{Z}}^0$ from \mathbf{Y}^0 and again obtain test statistics t_i^0 . Specifically, the test statistics t_i^0 for $i \in \mathcal{S}$ are saved as draws from the null distribution.

We repeat permutation procedure B times, and then utilize all saved sB permutation null statistics to calculate empirical p-values:

$$p_i = \frac{1}{sB} \sum_{b=1}^B \sum_{k \in \mathcal{S}_b} 1(t_k^{0b} \geq t_i).$$

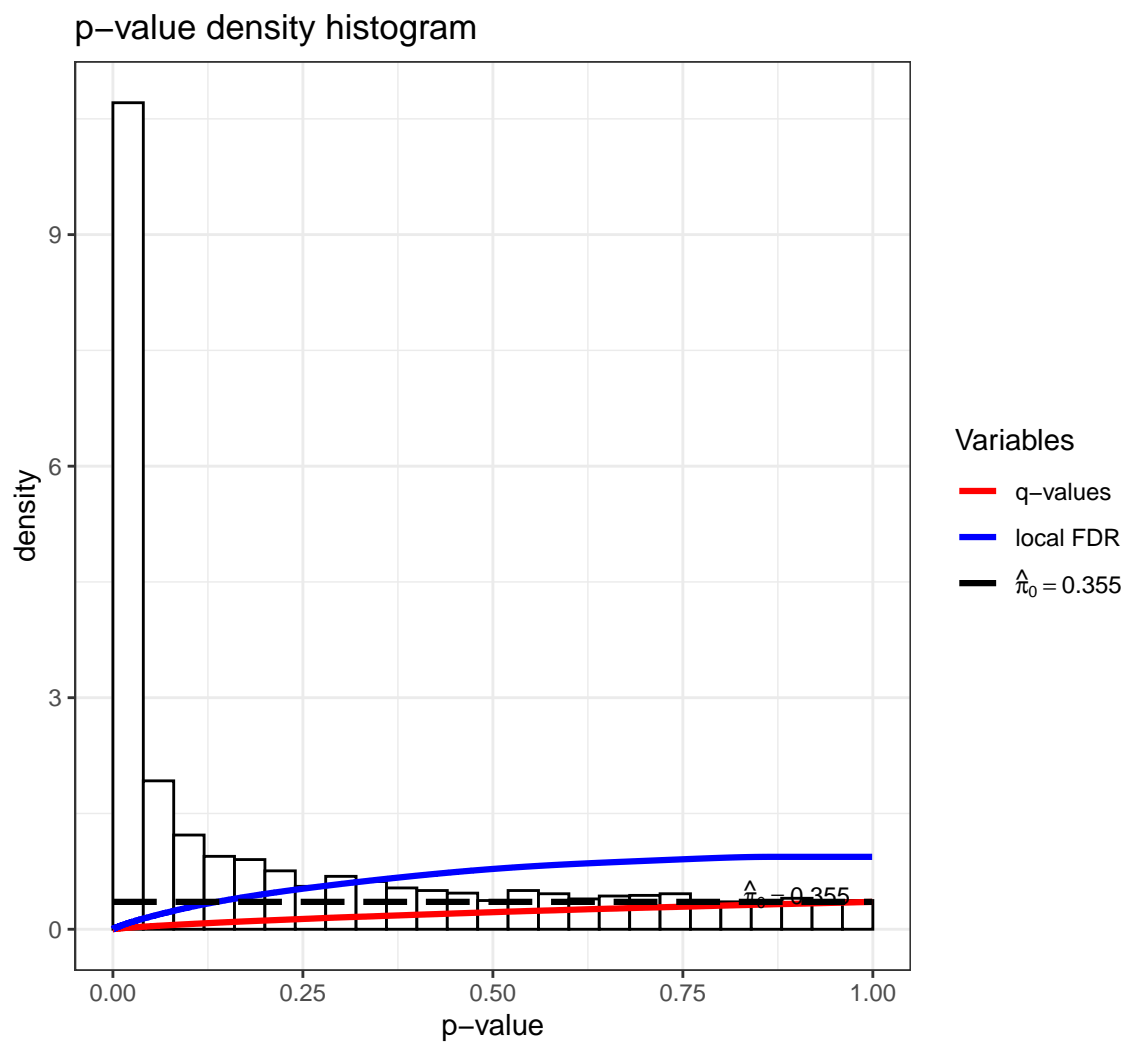
Example: Yeast Cell Cycle

Recall the yeast cell cycle data from earlier. We will test which genes have expression significantly associated with PC1 and PC2 since these both capture cell cycle regulation.

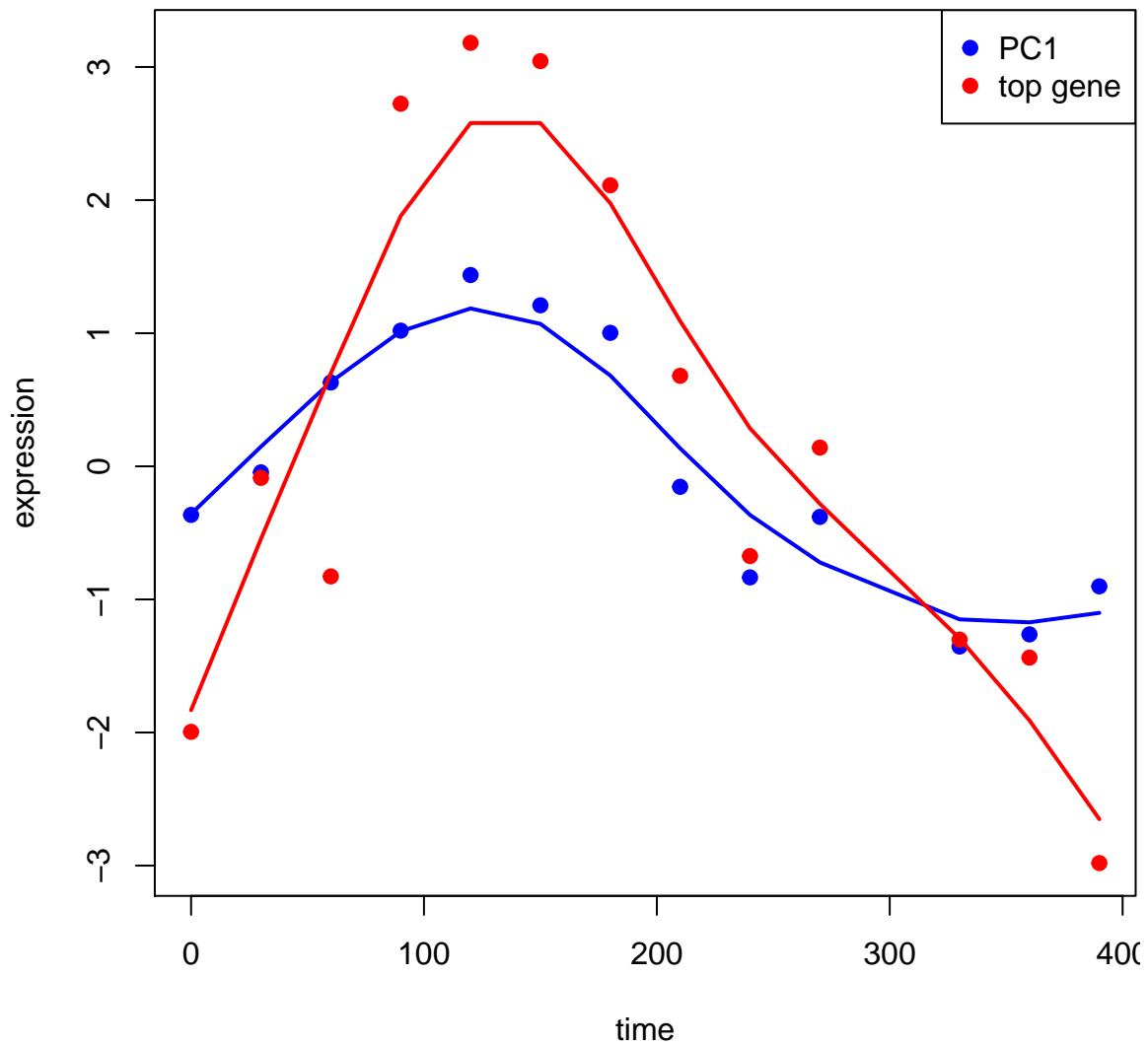
```
> library(jackstraw)
> load("../data/spellman.RData")
> time
[1] 0 30 60 90 120 150 180 210 240 270 330 360 390
> dim(gene_expression)
[1] 5981 13
> dat <- t(scale(t(gene_expression), center=TRUE, scale=FALSE))
```

Test for associations between PC1 and each gene, conditioning on PC1 and PC2 being relevant sources of systematic variation.

```
> jsobj <- jackstraw_pca(dat, r1=1, r=2, B=500, s=50, verbose=FALSE)
> jsobj$p.value %>% qvalue() %>% hist()
```

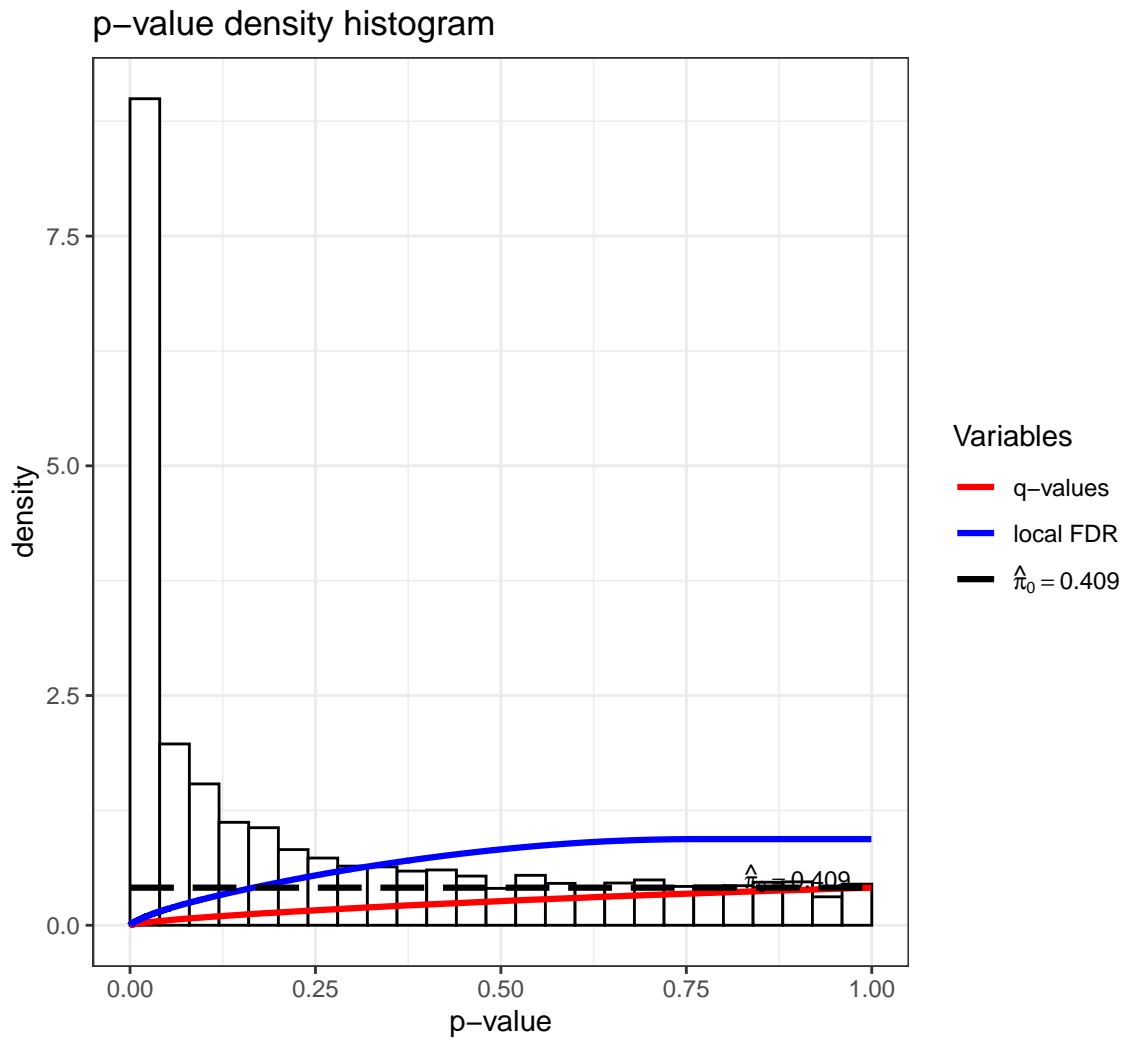


This is the most significant gene plotted with PC1.

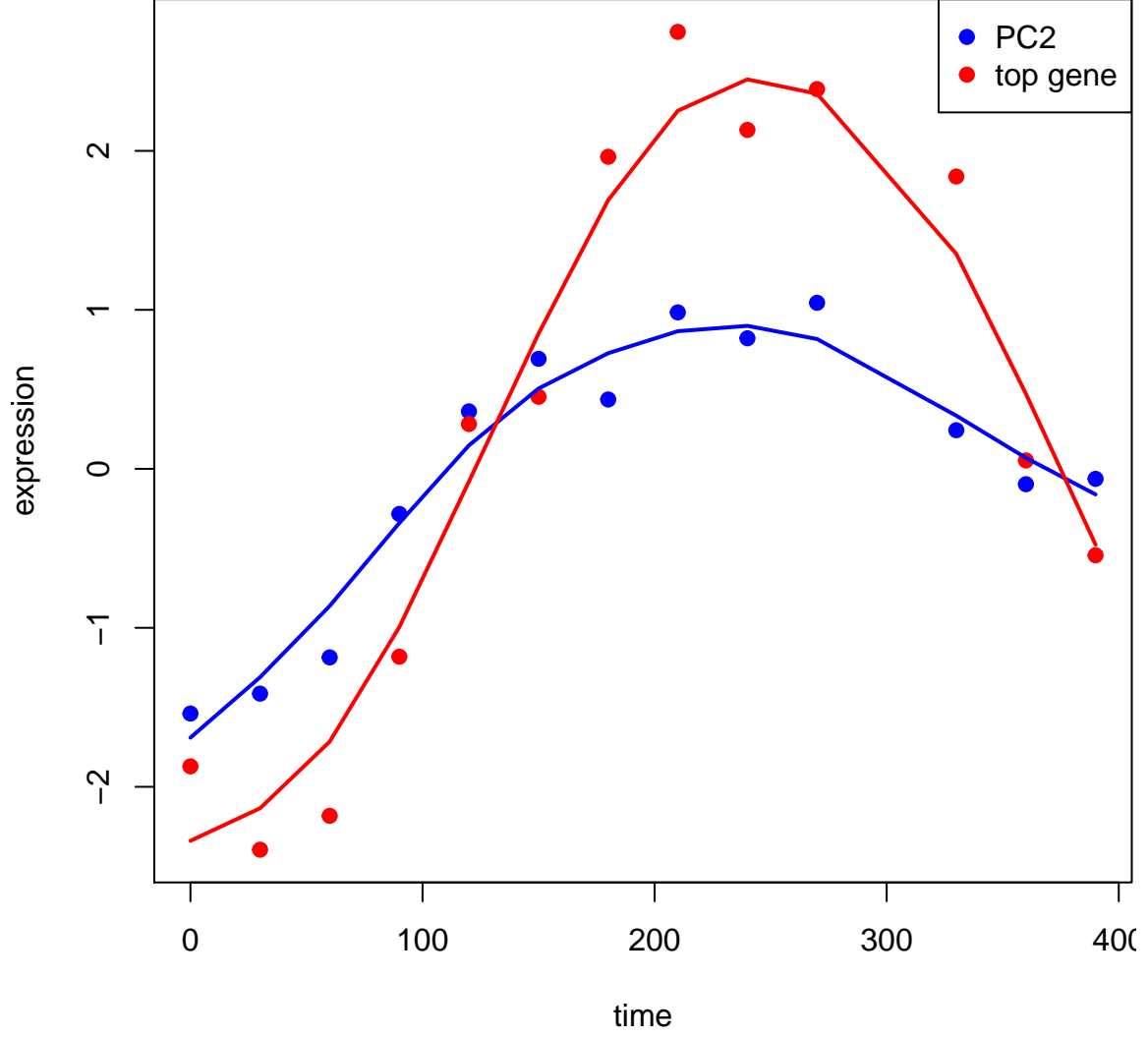


Test for associations between PC2 and each gene, conditioning on PC1 and PC2 being relevant sources of systematic variation.

```
> jsobj <- jackstraw_pca(dat, r1=2, r=2, B=500, s=50, verbose=FALSE)
> jsobj$p.value %>% qvalue() %>% hist()
```



This is the most significant gene plotted with PC2.



Surrogate Variable Analysis

The **surrogate variable analysis** (SVA) model combines the many responses model with the latent variable model introduced above:

$$\mathbf{Y}_{m \times n} = \mathbf{B}_{m \times d} \mathbf{X}_{d \times n} + \mathbf{\Phi}_{m \times r} \mathbf{Z}_{r \times n} + \mathbf{E}_{m \times n}$$

where $m \gg n > d + r$.

Here, only \mathbf{Y} and \mathbf{X} are observed, so we must combine many regressors model fitting techniques with latent variable estimation.

The variables \mathbf{Z} are called **surrogate variables** for what would be a complete model of all systematic variation.

Procedure

The main challenge is that the row spaces of \mathbf{X} and \mathbf{Z} may overlap. Even when \mathbf{X} is the result of a randomized experiment, there will be a high probability that the row spaces of \mathbf{X} and \mathbf{Z} have some overlap.

Therefore, one cannot simply estimate \mathbf{Z} by applying a latent variable estimation method on the residuals $\mathbf{Y} - \hat{\mathbf{B}}\mathbf{X}$ or on the observed response data \mathbf{Y} . In the former case, we will only estimate \mathbf{Z} in the space orthogonal to $\hat{\mathbf{B}}\mathbf{X}$. In the latter case, the estimate of \mathbf{Z} may modify the signal we can estimate in $\mathbf{B}\mathbf{X}$.

A recent method, takes an EM approach to estimating \mathbf{Z} in the model

$$\mathbf{Y}_{m \times n} = \mathbf{B}_{m \times d} \mathbf{X}_{d \times n} + \mathbf{\Phi}_{m \times r} \mathbf{Z}_{r \times n} + \mathbf{E}_{m \times n}.$$

It is shown to be necessary to penalize the likelihood in the estimation of \mathbf{B} — i.e., form shrinkage estimates of \mathbf{B} — in order to properly balance the row spaces of \mathbf{X} and \mathbf{Z} .

The regularized EM algorithm, called **cross-dimensonal inference** (CDI) iterates between

1. Estimate \mathbf{Z} from $\mathbf{Y} - \hat{\mathbf{B}}^{\text{Reg}} \mathbf{X}$
2. Estimate \mathbf{B} from $\mathbf{Y} - \hat{\mathbf{\Phi}} \hat{\mathbf{Z}}$

where $\hat{\mathbf{B}}^{\text{Reg}}$ is a regularized or shrunk estimate of \mathbf{B} .

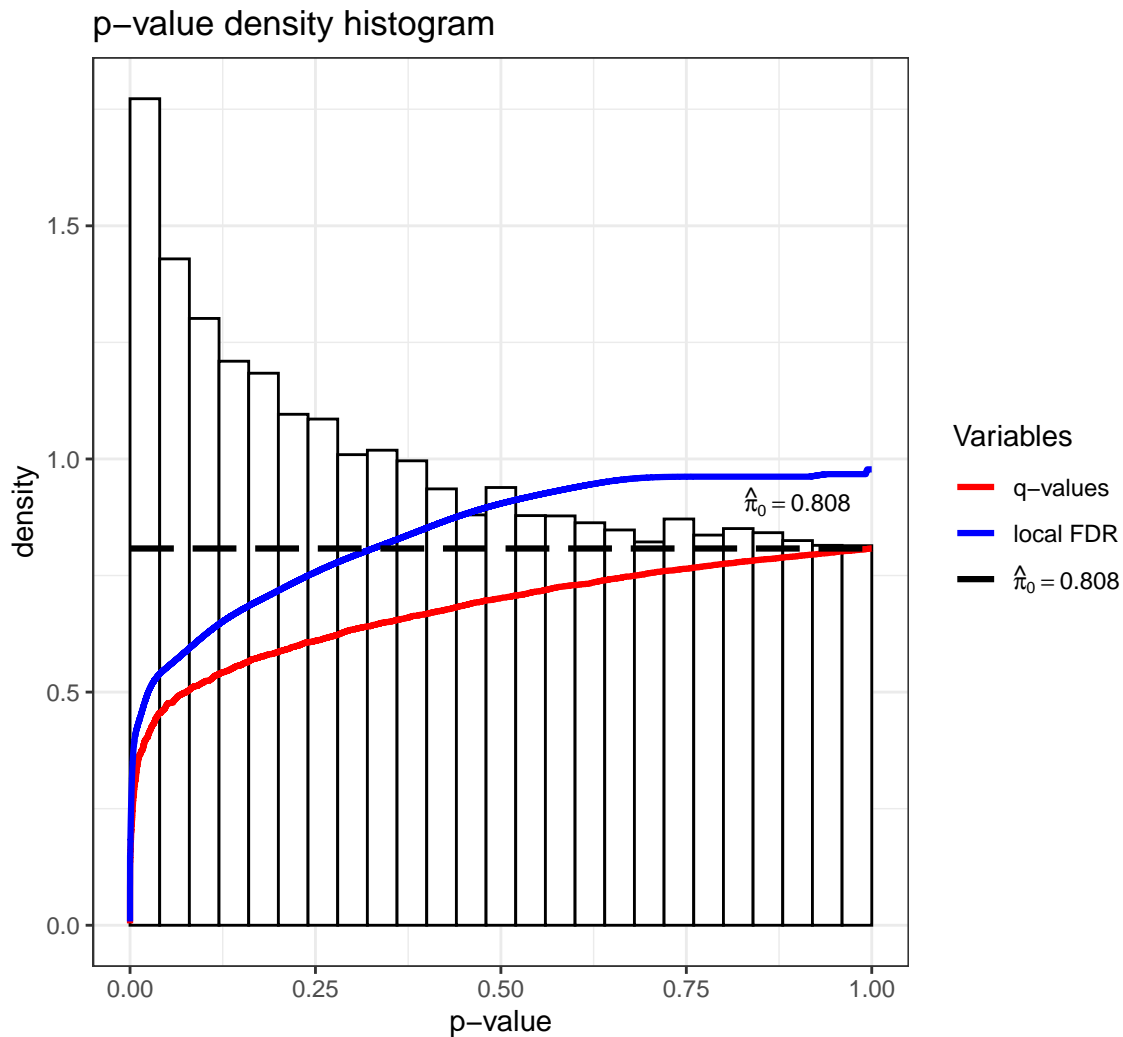
It can be shown that when the regularization can be represented by a prior distribution on \mathbf{B} then this algorithm achieves the MAP.

Example: Kidney Expr by Age

In Storey et al. (2005), we considered a study where kidney samples were obtained on individuals across a range of ages. The goal was to identify genes with expression associated with age.

```
> library(edge)
> library(splines)
> load("./data/kidney.RData")
> age <- kidcov$age
> sex <- kidcov$sex
> dim(kidexpr)
[1] 34061    72
> cov <- data.frame(sex = sex, age = age)
> null_model <- ~sex
> full_model <- ~sex + ns(age, df = 3)

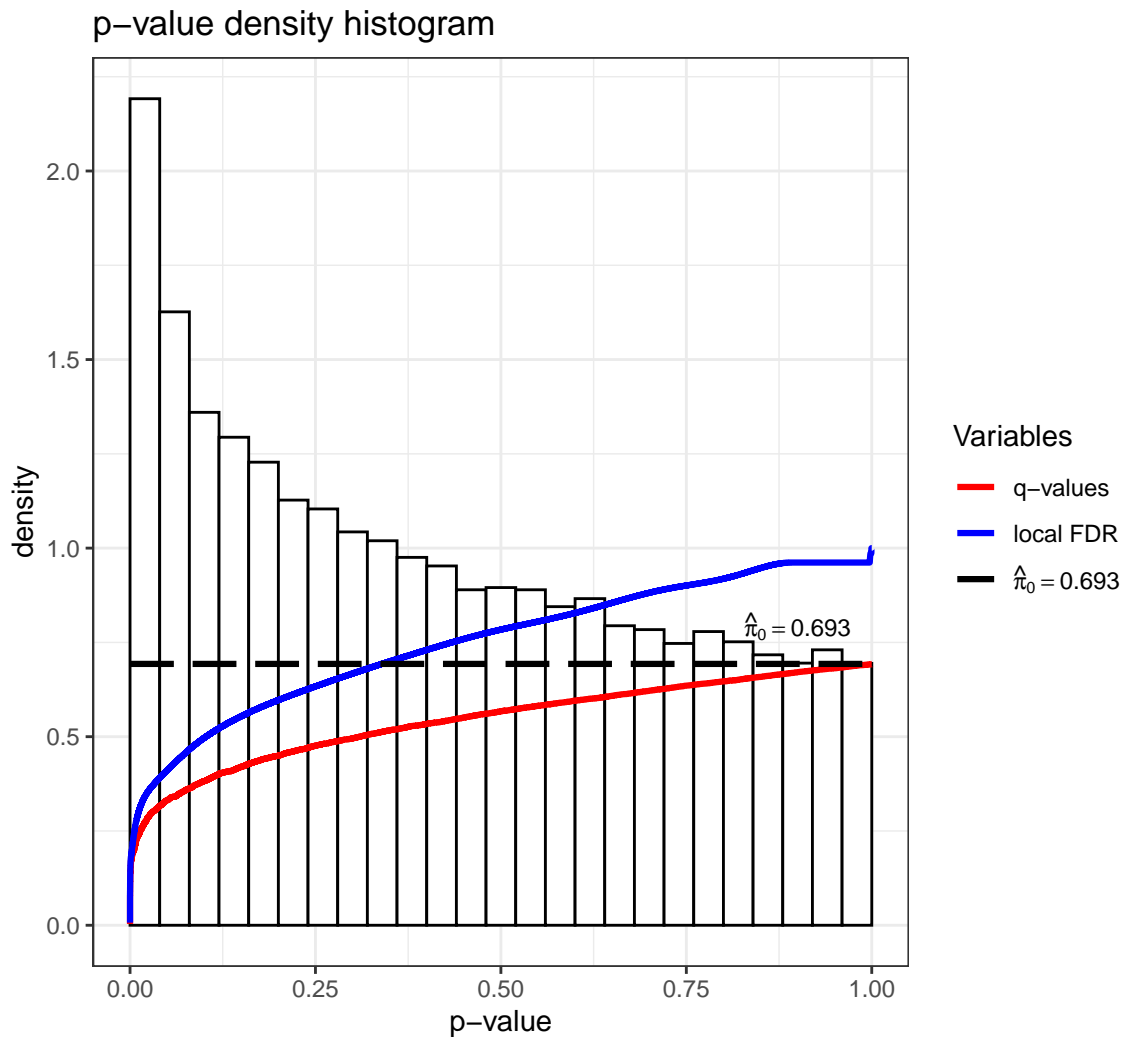
> de_obj <- build_models(data = kidexpr, cov = cov,
+                        null.model = null_model,
+                        full.model = full_model)
> de_lrt <- lrt(de_obj, nullDistn = "bootstrap", bs.its = 100, verbose=FALSE)
> qobj1 <- qvalueObj(de_lrt)
> hist(qobj1)
```



Now that we have completed a standard generalized LRT, let's estimate \mathbf{Z} (the surrogate variables) using the `sva` package as accessed via the `edge` package.

```
> dim(nullMatrix(de_obj))
[1] 72 2
> de_sva <- apply_sva(de_obj, n.sv=4, method="irw", B=10)
Number of significant surrogate variables is: 4
Iteration (out of 10): 1 2 3 4 5 6 7 8 9 10
> dim(nullMatrix(de_sva))
[1] 72 6
> de_svalrt <- lrt(de_sva, nullDistn = "bootstrap", bs.its = 100, verbose=FALSE)

> qobj2 <- qvalueObj(de_svalrt)
> hist(qobj2)
```

```
> summary(qobj1)

Call:
qvalue(p = pval)

pi0:    0.8081212

Cumulative number of significant calls:

      <1e-04 <0.001 <0.01 <0.025 <0.05 <0.1    <1
p-value    27    161    798   1676   2906  5271 34061
q-value     0     0     2     4     10   27 34061
local FDR   0     0     2     2     5   18 34061
```

```
> summary(qobj2)

Call:
qvalue(p = pval)

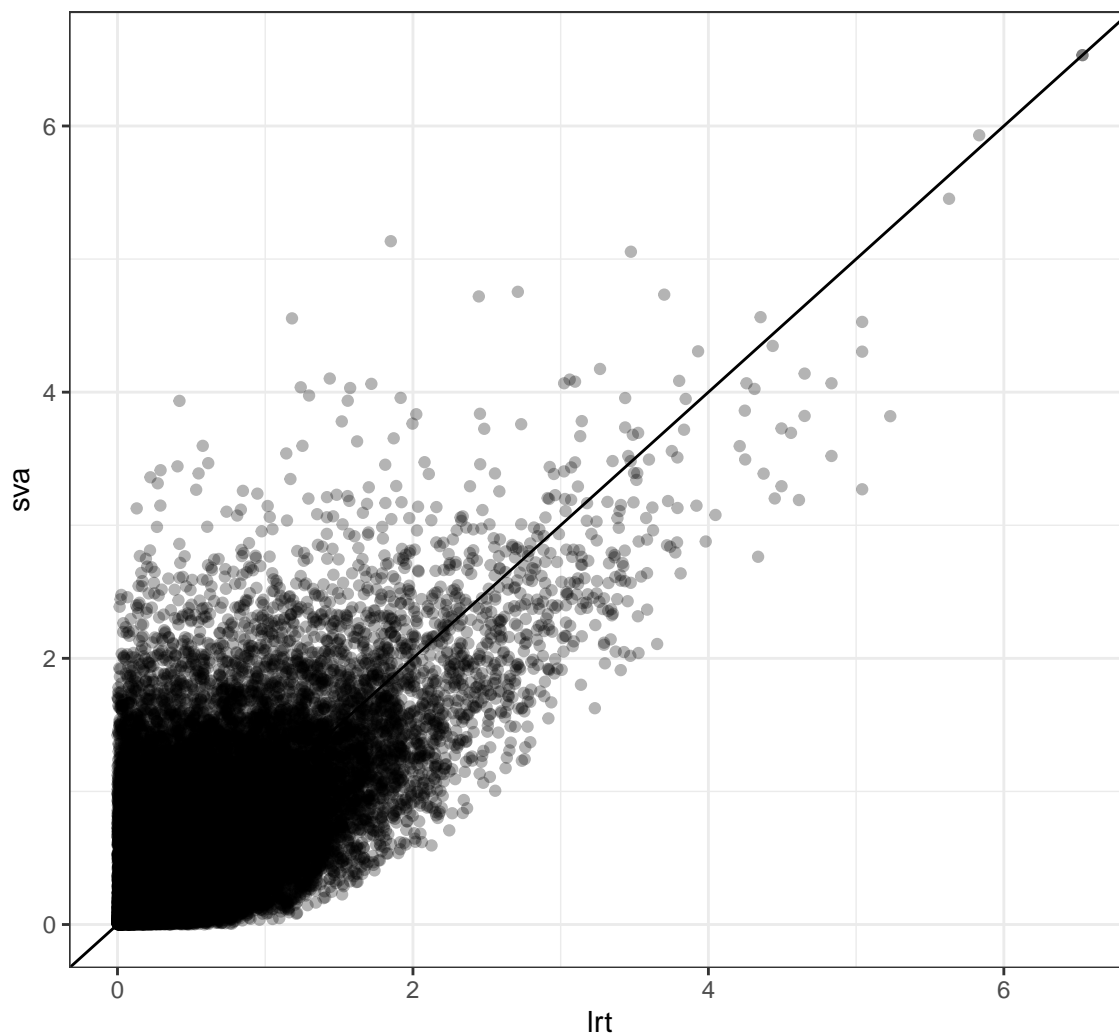
pi0:    0.6925105
```

Cumulative number of significant calls:

	<1e-04	<0.001	<0.01	<0.025	<0.05	<0.1	<1
p-value	28	151	1001	2051	3549	6168	34061
q-value	0	0	3	4	6	51	34061
local FDR	0	0	2	2	3	28	34053

P-values from two analyses are fairly different.

```
> data.frame(lrt=-log10(qobj1$pval), sva=-log10(qobj2$pval)) %>%  
+   ggplot() + geom_point(aes(x=lrt, y=sva), alpha=0.3) + geom_abline()
```



Extras

Source

License

Source Code

Session Information

```
> sessionInfo()
R version 3.6.0 (2019-04-26)
Platform: x86_64-apple-darwin15.6.0 (64-bit)
Running under: macOS 10.15.3

Matrix products: default
BLAS: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRblas.0.dylib
LAPACK: /Library/Frameworks/R.framework/Versions/3.6/Resources/lib/libRlapack.dylib

locale:
[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8

attached base packages:
[1] splines parallel stats graphics grDevices utils
[7] datasets methods base

other attached packages:
[1] edge_2.14.0 Biobase_2.42.0
[3] BiocGenerics_0.28.0 jackstraw_1.3
[5] qvalue_2.15.0 MASS_7.3-51.5
[7] broom_0.5.2 forcats_0.5.0
[9] stringr_1.4.0 dplyr_0.8.4
[11] purrr_0.3.3 readr_1.3.1
[13] tidyr_1.0.2 tibble_2.1.3
[15] ggplot2_3.2.1 tidyverse_1.3.0
[17] knitr_1.28

loaded via a namespace (and not attached):
[1] nlme_3.1-144 matrixStats_0.55.0
[3] bitops_1.0-6 fs_1.3.1
[5] bit64_0.9-7 lubridate_1.7.4
[7] httr_1.4.1 tools_3.6.0
[9] backports_1.1.5 R6_2.4.1
[11] irlba_2.3.3 mgcv_1.8-31
[13] DBI_1.1.0 lazyeval_0.2.2
[15] colorspace_1.4-1 withr_2.1.2
[17] tidyselect_1.0.0 bit_1.1-15.2
[19] compiler_3.6.0 cli_2.0.2
[21] rvest_0.3.5 xml2_1.2.2
[23] labeling_0.3 scales_1.1.0
[25] genefilter_1.64.0 digest_0.6.25
[27] minqa_1.2.4 rmarkdown_2.1
[29] pkgconfig_2.0.3 htmltools_0.4.0
[31] lme4_1.1-21 limma_3.38.3
[33] dbplyr_1.4.2 rlang_0.4.5
[35] readxl_1.3.1 RSQLite_2.2.0
[37] rstudioapi_0.11 farver_2.0.3
[39] generics_0.0.2 jsonlite_1.6.1
[41] BiocParallel_1.16.5 gtools_3.8.1
[43] RCurl_1.98-1.1 magrittr_1.5
[45] Matrix_1.2-18 Rcpp_1.0.3
[47] munsell_0.5.0 S4Vectors_0.20.1
```

[49]	fansi_0.4.1	lifecycle_0.1.0
[51]	stringi_1.4.6	yaml_2.2.1
[53]	ClusterR_1.2.1	plyr_1.8.5
[55]	blob_1.2.1	grid_3.6.0
[57]	crayon_1.3.4	lattice_0.20-40
[59]	haven_2.2.0	annotate_1.60.0
[61]	hms_0.5.3	pillar_1.4.3
[63]	boot_1.3-24	corpcor_1.6.9
[65]	codetools_0.2-16	stats4_3.6.0
[67]	reshape2_1.4.3	XML_3.99-0.3
[69]	reprex_0.3.0	glue_1.3.1
[71]	evaluate_0.14	snm_1.30.0
[73]	modelr_0.1.6	vctrs_0.2.3
[75]	nloptr_1.2.1	cellranger_1.1.0
[77]	gtable_0.3.0	assertthat_0.2.1
[79]	xfun_0.12	rsvd_1.0.3
[81]	lfa_1.12.0	xtable_1.8-4
[83]	survival_3.1-8	IRanges_2.16.0
[85]	memoise_1.1.0	AnnotationDbi_1.44.0
[87]	cluster_2.1.0	sva_3.30.1
[89]	gmp_0.5-13.6	