

Regularized Precision Matrix Estimation via ADMM

Matt Galloway

April 1, 2018

Abstract

ADMMsigma is an R package that estimates a penalized precision matrix via the alternating direction method of multipliers (ADMM) algorithm. It currently supports a general elastic-net penalty that allows for both ridge and lasso-type penalties as special cases. This report will provide a brief overview of the algorithm and detail how it can be utilized to estimate precision matrices of jointly normal distributions. In addition, examples and simulation results will be provided.

Contents

1	Introduction	1
2	Regularized Precision Matrix Estimation	2
2.1	Condensed-Form ADMM	3
2.2	Algorithm	4
3	R Package	6
3.1	Installation	6
3.2	Usage	7
3.3	Benchmark	11
3.4	Simulations	13
	References	19

1 Introduction

Suppose we want to solve the following optimization problem:

$$\begin{aligned} & \text{minimize } f(x) + g(z) \\ & \text{subject to } Ax + Bz = c \end{aligned}$$

where $x \in \mathbb{R}^n$, $z \in \mathbb{R}^m$, $A \in \mathbb{R}^{p \times n}$, $B \in \mathbb{R}^{p \times m}$, $c \in \mathbb{R}^p$ – though we will later consider cases where x and z are matrices. Further, we will assume f and g are convex. The *augmented lagrangian* is constructed as follows:

$$L_\rho(x, z, y) = f(x) + g(z) + y^T(Ax + Bz - c) + \frac{\rho}{2} \|Ax + Bz - c\|_2^2$$

where $y \in \mathbb{R}^p$ is the lagrange multiplier. The optimal value is

$$p^* = \inf \{f(x) + g(z) | Ax + Bz = c\}$$

Clearly, the minimization under the augmented lagrangian is equivalent to that of the usual lagrangian since any feasible point (x, z) satisfies the constraint $\rho \|Ax + Bz - c\|_2^2 / 2 = 0$.

The alternating direct method of multipliers (ADMM) algorithm consists of the following repeated iterations:

$$x^{k+1} := \arg \min_x L_\rho(x, z^k, y^k) \quad (1)$$

$$z^{k+1} := \arg \min_z L_\rho(z^{k+1}, z, y^k) \quad (2)$$

$$y^{k+1} := y^k + \rho(Ax^{k+1} + Bz^{k+1} - c) \quad (3)$$

A more complete introduction to the algorithm – specifically how it arose out of *dual ascent* and *method of multipliers* – can be found in Boyd et al. (2011).

2 Regularized Precision Matrix Estimation

We now consider the case where X_1, \dots, X_n are iid $N_p(\mu, \Sigma)$ and we are tasked with estimating the precision matrix, denoted $\Omega \equiv \Sigma^{-1}$. The maximum likelihood estimator for Ω is

$$\hat{\Omega} = \arg \min_{\Omega \in S_+^p} \{Tr(S\Omega) - \log \det(\Omega)\}$$

where $S = \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})^T / n$. It is straight forward to show that when the solution exists, $\hat{\Omega} = S^{-1}$.

We can construct a *penalized* likelihood estimator by adding a penalty term, $P(\Omega)$, to the likelihood:

$$\hat{\Omega}_\lambda = \arg \min_{\Omega \in S_+^p} \{Tr(S\Omega) - \log \det(\Omega) + P(\Omega)\}$$

Throughout the rest of this document we will take $P(\Omega)$ to be $P(\Omega) = \lambda \left[\frac{1-\alpha}{2} \|\Omega\|_F^2 + \alpha \|\Omega\|_1 \right]$ so that the full penalized likelihood is as follows:

$$\hat{\Omega}_\lambda = \arg \min_{\Omega \in S_+^p} \left\{ Tr(S\Omega) - \log \det(\Omega) + \lambda \left[\frac{1-\alpha}{2} \|\Omega\|_F^2 + \alpha \|\Omega\|_1 \right] \right\}$$

where $0 \leq \alpha \leq 1$, $\lambda > 0$, $\|\cdot\|_F^2$ is the Frobenius norm and we define $\|A\|_1 = \sum_{i,j} |A_{ij}|$. This *elastic-net* penalty was explored by Hui Zou and Trevor Hastie (Zou and Hastie 2005) and is identical to the penalty used in the popular penalized regression package **glmnet**. Clearly, when $\alpha = 0$ the elastic-net reduces to a ridge-type penalty and when $\alpha = 1$ this reduces to a lasso-type penalty.

By letting f be equal to the non-penalized likelihood and g equal to $P(\Omega)$, our goal is to minimize the full augmented lagrangian where the constraint is that $\Omega - Z$ is equal to zero:

$$L_\rho(\Omega, Z, \Lambda) = f(\Omega) + g(Z) + Tr[\Lambda(\Omega - Z)] + \frac{\rho}{2} \|\Omega - Z\|_F^2$$

The ADMM algorithm for regularized precision matrix estimation is

$$\Omega^{k+1} = \arg \min_{\Omega} \left\{ \text{Tr}(\Omega) - \log \det(\Omega) + \text{Tr}[\Lambda^k (\Omega - Z^k)] + \frac{\rho}{2} \|\Omega - Z^k\|_F^2 \right\} \quad (4)$$

$$Z^{k+1} = \arg \min_Z \left\{ \lambda \left[\frac{1-\alpha}{2} \|Z\|_F^2 + \alpha \|Z\|_1 \right] + \text{Tr}[\Lambda^k (\Omega^{k+1} - Z)] + \frac{\rho}{2} \|\Omega^{k+1} - Z\|_F^2 \right\} \quad (5)$$

$$\Lambda^{k+1} = \Lambda^k + \rho (\Omega^{k+1} - Z^{k+1}) \quad (6)$$

2.1 Condensed-Form ADMM

An alternate form of the ADMM algorithm can be constructed by scaling the dual variable. Let us define $R^k = \Omega - Z^k$ and $U^k = \Lambda^k / \rho$. Then

$$\begin{aligned} \text{Tr}[\Lambda^k (\Omega - Z^k)] + \frac{\rho}{2} \|\Omega - Z^k\|_F^2 &= \text{Tr}[\Lambda^k R^k] + \frac{\rho}{2} \|R^k\|_F^2 \\ &= \frac{\rho}{2} \|R^k + \Lambda^k / \rho\|_F^2 - \frac{\rho}{2} \|\Lambda^k / \rho\|_F^2 \\ &= \frac{\rho}{2} \|R^k + U^k\|_F^2 - \frac{\rho}{2} \|U^k\|_F^2 \end{aligned}$$

The condensed-form can now be written as follows:

$$\Omega^{k+1} = \arg \min_{\Omega} \left\{ \text{Tr}(\Omega) - \log \det(\Omega) + \frac{\rho}{2} \|\Omega - Z^k + U^k\|_F^2 \right\} \quad (7)$$

$$Z^{k+1} = \arg \min_Z \left\{ \lambda \left[\frac{1-\alpha}{2} \|Z\|_F^2 + \alpha \|Z\|_1 \right] + \frac{\rho}{2} \|\Omega^{k+1} - Z + U^k\|_F^2 \right\} \quad (8)$$

$$U^{k+1} = U^k + \Omega^{k+1} - Z^{k+1} \quad (9)$$

More generally (in vector form),

$$x^{k+1} := \arg \min_x \left\{ f(x) + \frac{\rho}{2} \|Ax + Bz^k - c + u^k\|_2^2 \right\} \quad (10)$$

$$z^{k+1} := \arg \min_z \left\{ g(z) + \frac{\rho}{2} \|Ax^{k+1} + Bz - c + u^k\|_2^2 \right\} \quad (11)$$

$$u^{k+1} := u^k + Ax^{k+1} + Bz^{k+1} - c \quad (12)$$

Note that there are limitations to using this method. For instance, because the dual variable is scaled by ρ (the step size), this form limits one to using a constant step size (without making further adjustments to U^k) – a limitation that could prolong the convergence rate. Because of this, we will only consider the non-condensed form for the remainder of this report.

2.2 Algorithm

$$\begin{aligned}\Omega^{k+1} &= \arg \min_{\Omega} \left\{ Tr(\Omega) - \log \det(\Omega) + Tr[\Lambda^k (\Omega - Z^k)] + \frac{\rho}{2} \|\Omega - Z^k\|_F^2 \right\} \\ Z^{k+1} &= \arg \min_Z \left\{ \lambda \left[\frac{1-\alpha}{2} \|Z\|_F^2 + \alpha \|Z\|_1 \right] + Tr[\Lambda^k (\Omega^{k+1} - Z)] + \frac{\rho}{2} \|\Omega^{k+1} - Z\|_F^2 \right\} \\ \Lambda^{k+1} &= \Lambda^k + \rho (\Omega^{k+1} - Z^{k+1})\end{aligned}$$

1. Decompose $S + \Lambda^k - \rho Z^k = VQV^T$.

$$\Omega^{k+1} = \frac{1}{2\rho} V \left[-Q + (Q^2 + 4\rho I_p)^{1/2} \right] V^T$$

2. Elementwise soft-thresholding for all $i = 1, \dots, p$ and $j = 1, \dots, p$.

$$\begin{aligned}Z_{ij}^{k+1} &= \frac{1}{\lambda(1-\alpha) + \rho} \text{sign}(\rho\Omega_{ij}^{k+1} + \Lambda_{ij}^k) (|\rho\Omega_{ij}^{k+1} + \Lambda_{ij}^k| - \lambda\alpha)_+ \\ &= \frac{1}{\lambda(1-\alpha) + \rho} \text{Soft}((\rho\Omega_{ij}^{k+1} + \Lambda_{ij}^k), \lambda\alpha)\end{aligned}$$

3. Update Λ .

$$\Lambda^{k+1} = \Lambda^k + \rho (\Omega^{k+1} - Z^{k+1})$$

2.2.1 Proof of (1):

$$\Omega^{k+1} = \arg \min_{\Omega} \left\{ Tr(\Omega) - \log \det(\Omega) + Tr[\Lambda^k (\Omega - Z^k)] + \frac{\rho}{2} \|\Omega - Z^k\|_F^2 \right\}$$

Code snippet:

Note this is not the actual code. The real code is written in c++.

```
# ridge penalized precision matrix
# function
RIDGEsigma = function(S, lam) {

  # dimensions
  p = dim(S)[1]

  # gather eigen values of S (spectral
  # decomposition)
  e.out = eigen(S, symmetric = TRUE)
```

```

# augment eigen values for omega hat
new.evs = (-e.out$val + sqrt(e.out$val^2 +
  4 * lam))/(2 * lam)

# compute omega hat for lambda (zero
# gradient equation)
omega = tcrossprod(e.out$vec * rep(new.evs,
  each = p), e.out$vec)

# compute gradient
grad = S - qr.solve(omega) + lam * omega

return(list(omega = omega, gradient = grad))
}

```

2.2.2 Proof of (2)

$$Z^{k+1} = \arg \min_Z \left\{ \lambda \left[\frac{1-\alpha}{2} \|Z\|_F^2 + \alpha \|Z\|_1 \right] + Tr [\Lambda^k (\Omega^{k+1} - Z)] + \frac{\rho}{2} \|\Omega^{k+1} - Z\|_F^2 \right\}$$

Code snippet:

Note this is not the actual code. The real code is written in c++.

```

# ADMMsigma function
ADMMsigma = function(X = NULL, S = NULL,
  lam, alpha = 1, rho = 2, mu = 10, tau1 = 2,
  tau2 = 2, tol1 = 1e-04, tol2 = 1e-04,
  maxit = 1000) {

  # compute sample covariance matrix, if
  # necessary
  if (is.null(S)) {

    # covariance matrix
    n = dim(X)[1]
    S = (n - 1)/n * cov(X)

  }

  # allocate memory
  p = dim(S)[1]
  criterion = TRUE
  iter = lik = s = r = eps1 = eps2 = 0
  new.Z = Y = Omega = matrix(0, nrow = p,
    ncol = p)

  # loop until convergence

```

```

while (criterion && (iter <= maxit)) {

  # ridge equation (1) gather eigen values
  # (spectral decomposition)
  Z = new.Z
  Omega = sigma_ridge(S + Y - rho *
    Z, lam = rho)$omega

  # penalty equation (2) soft-thresholding
  new.Z = soft(Y + rho * Omega, lam *
    alpha)/(lam * (1 - alpha) + rho)

  # update U (3)
  Y = Y + rho * (Omega - new.Z)

  # calculate new rho
  s = sqrt(sum((rho * (new.Z - Z))^2))
  r = sqrt(sum((Omega - new.Z)^2))
  rho = rho * (tau1 * (r > mu * s) +
    (s > mu * r)/tau2 + (s/mu <=
    r & r <= mu * s))
  iter = iter + 1

  # stopping criterion
  eps1 = p * tol1 + tol2 * max(sqrt(sum(Omega^2)),
    sqrt(sum(new.Z^2)))
  eps2 = p * tol1 + tol2 * sqrt(sum(Y^2))
  criterion = (r >= eps1 || s >= eps2)

}
return(list(Iterations = iter, Omega = Omega))
}

```

3 R Package

3.1 Installation

```

# The easiest way to install is from CRAN
install.packages("ADMMsigma")

# You can also install the development
# version from GitHub:
# install.packages('devtools')
devtools::install_github("MGallow/ADMMsigma")

```

If there are any issues/bugs, please let me know: [github](#). You can also contact me via my website. Pull requests are welcome!

A (possibly incomplete) list of functions contained in the package can be found below:

- `ADMMsigma()` computes the estimated precision matrix (ridge, lasso, and elastic-net type regularization optional)
- `RIDGEsigma()` computes the estimated ridge penalized precision matrix via closed-form solution
- `plot.ADMMsigma()` produces a heat map for cross validation errors
- `plot.RIDGEsigma()` produces a heat map for cross validation errors

3.2 Usage

```
library(ADMMsigma)

# generate data from a sparse matrix
# first compute covariance matrix
S = matrix(0.7, nrow = 5, ncol = 5)
for (i in 1:5) {
  for (j in 1:5) {
    S[i, j] = S[i, j]^abs(i - j)
  }
}

# print oracle precision matrix
# (shrinkage might be useful)
(Omega = qr.solve(S) %>% round(3))

##          [,1] [,2] [,3] [,4] [,5]
## [1,]  1.961 -1.373  0.000  0.000  0.000
## [2,] -1.373  2.922 -1.373  0.000  0.000
## [3,]  0.000 -1.373  2.922 -1.373  0.000
## [4,]  0.000  0.000 -1.373  2.922 -1.373
## [5,]  0.000  0.000  0.000 -1.373  1.961

# generate 1000 x 5 matrix with rows
# drawn from iid  $N_p(0, S)$ 
Z = matrix(rnorm(1000 * 5), nrow = 1000,
           ncol = 5)
out = eigen(S, symmetric = TRUE)
S.sqrt = out$vectors %*% diag(out$values^0.5) %*%
         t(out$vectors)
X = Z %*% S.sqrt

# print sample precision matrix (perhaps
# a bad estimate)
(qr.solve(cov(X)) %>% round(5))
```

```
##          [,1] [,2] [,3] [,4] [,5]
```

```
## [1,] 1.99991 -1.41471 0.09295 -0.08643 0.06678
## [2,] -1.41471 2.98077 -1.39085 -0.00180 -0.02666
## [3,] 0.09295 -1.39085 2.93466 -1.57182 0.13492
## [4,] -0.08643 -0.00180 -1.57182 3.28802 -1.50476
## [5,] 0.06678 -0.02666 0.13492 -1.50476 1.95149
```

```
# elastic-net type penalty (set tolerance
# to 1e-8)
ADMMsigma(X, tol1 = 1e-08, tol2 = 1e-08)
```

```
##
## Iterations:
## [1] 38
##
## Tuning parameters:
##      log10(lam)  alpha
## [1,]         -3    0.6
##
## Omega:
##      [,1]      [,2]      [,3]      [,4]      [,5]
## [1,] 1.98552 -1.39091 0.07533 -0.07544 0.06043
## [2,] -1.39091 2.94059 -1.35892 -0.02041 -0.01718
## [3,] 0.07533 -1.35892 2.89522 -1.53331 0.11552
## [4,] -0.07544 -0.02041 -1.53331 3.23632 -1.47668
## [5,] 0.06043 -0.01718 0.11552 -1.47668 1.93602
```

```
# lasso penalty (default tolerance)
ADMMsigma(X, alpha = 1)
```

```
##
## Iterations:
## [1] 16
##
## Tuning parameters:
##      log10(lam)  alpha
## [1,]         -3     1
##
## Omega:
##      [,1]      [,2]      [,3]      [,4]      [,5]
## [1,] 1.98897 -1.39520 0.07501 -0.07195 0.05755
## [2,] -1.39520 2.94945 -1.36518 -0.02047 -0.01395
## [3,] 0.07501 -1.36518 2.90753 -1.54178 0.11625
## [4,] -0.07195 -0.02047 -1.54178 3.24849 -1.48219
## [5,] 0.05755 -0.01395 0.11625 -1.48219 1.93985
```

```
# elastic-net penalty (alpha = 0.5)
ADMMsigma(X, alpha = 0.5)
```

```
##
## Iterations:
## [1] 16
##
## Tuning parameters:
##      log10(lam)  alpha
## [1,]         -3    0.5
##
```



```
## Omega:
##      [,1]      [,2]      [,3]      [,4]      [,5]
## [1,]  1.98401 -1.38850  0.07368 -0.07468  0.06041
## [2,] -1.38850  2.93561 -1.35372 -0.02388 -0.01640
## [3,]  0.07368 -1.35372  2.88711 -1.52619  0.11302
## [4,] -0.07468 -0.02388 -1.52619  3.22819 -1.47294
## [5,]  0.06041 -0.01640  0.11302 -1.47294  1.93396
```

```
# ridge penalty
ADMMsigma(X, alpha = 0)
```

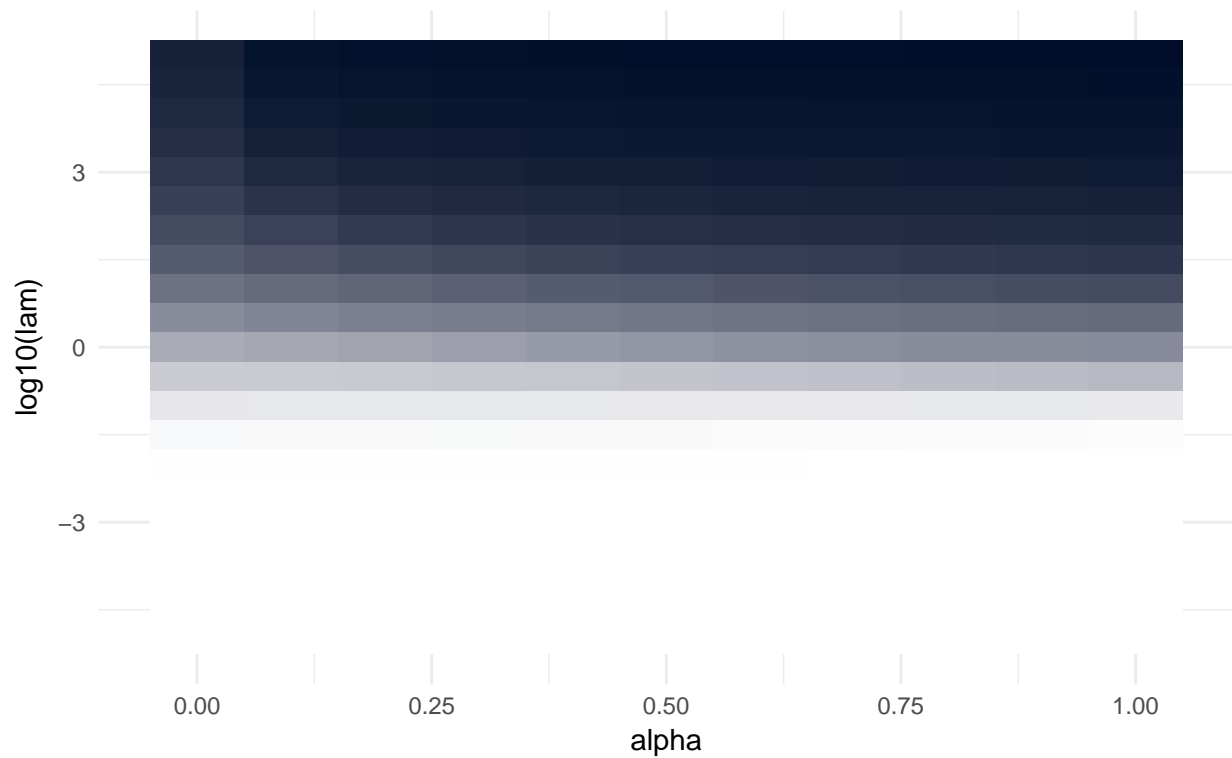
```
##
## Iterations:
## [1] 16
##
## Tuning parameters:
##      log10(lam)  alpha
## [1,]         -3      0
##
## Omega:
##      [,1]      [,2]      [,3]      [,4]      [,5]
## [1,]  1.97957 -1.38283  0.07351 -0.07837  0.06365
## [2,] -1.38283  2.92392 -1.34491 -0.02497 -0.01979
## [3,]  0.07351 -1.34491  2.87050 -1.51430  0.11142
## [4,] -0.07837 -0.02497 -1.51430  3.21177 -1.46548
## [5,]  0.06365 -0.01979  0.11142 -1.46548  1.92892
```

```
# ridge penalty no ADMM
RIDGEsigma(X, lam = 10^seq(-8, 8, 0.01))
```

```
##
## Tuning parameter:
##      log10(lam)    lam
## [1,]        -2.82  0.002
##
## Omega:
##      [,1]      [,2]      [,3]      [,4]      [,5]
## [1,]  1.96954 -1.36852  0.06659 -0.07680  0.06314
## [2,] -1.36852  2.89875 -1.32718 -0.03111 -0.01877
## [3,]  0.06659 -1.32718  2.84508 -1.49281  0.10334
## [4,] -0.07680 -0.03111 -1.49281  3.18037 -1.44902
## [5,]  0.06314 -0.01877  0.10334 -1.44902  1.91851
```

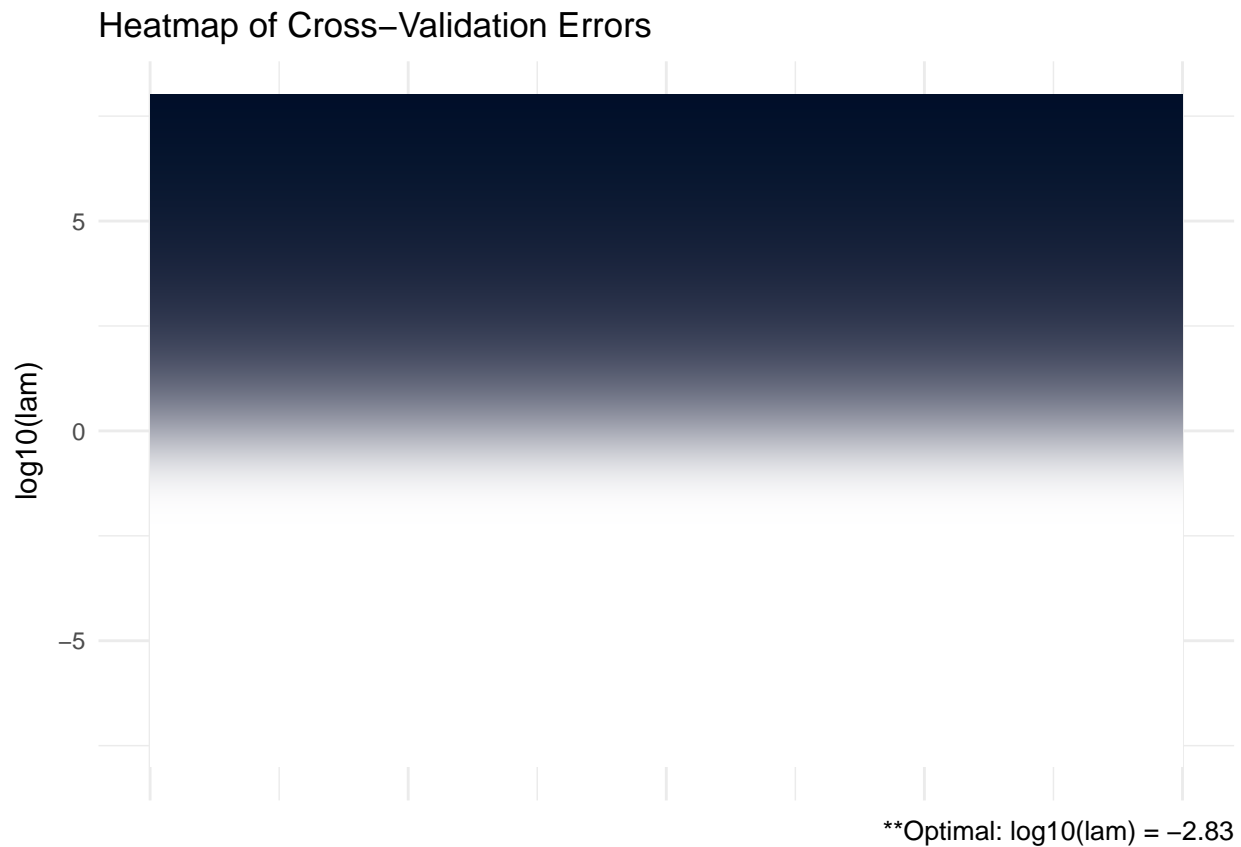
```
# produce CV heat map for ADMMsigma
ADMMsigma(X, tol1 = 1e-08, tol2 = 1e-08) %>%
  plot
```

Heatmap of Cross-Validation Errors



****Optimal: $\log_{10}(\text{lam}) = -3$, $\alpha = 0.8$**

```
# produce CV heat map for RIDGESigma
RIDGESigma(X, lam = 10^seq(-8, 8, 0.01)) %>%
  plot
```



3.3 Benchmark

3.3.1 Computer Specs:

- MacBook Pro (Late 2016)
- Processor: 2.9 GHz Intel Core i5
- Memory: 8GB 2133 MHz
- Graphics: Intel Iris Graphics 550

```
# generate data from tri-diagonal  
# (sparse) matrix compute covariance  
# matrix (can confirm inverse is  
# tri-diagonal)  
S = matrix(0, nrow = 100, ncol = 100)  
  
for (i in 1:100) {  
  for (j in 1:100) {  
    S[i, j] = 0.7^(abs(i - j))  
  }  
}  
  
# generate 1000 x 100 matrix with rows  
# drawn from iid N_p(0, S)
```

```

Z = matrix(rnorm(1000 * 100), nrow = 1000,
           ncol = 100)
out = eigen(S, symmetric = TRUE)
S.sqrt = out$vectors %*% diag(out$values^0.5) %*%
         t(out$vectors)
X = Z %*% S.sqrt

# glasso (for comparison)
microbenchmark(glasso(s = S, rho = 0.1))

## Unit: milliseconds
##           expr      min       lq      mean    median      uq
##  glasso(s = S, rho = 0.1) 49.49746 51.2968 55.21674 53.17891 56.54074
##           max neval
##  94.10694   100

# benchmark ADMMsigma - default tolerance
microbenchmark(ADMMsigma(S = S, lam = 0.1,
                        alpha = 1, tol1 = 1e-04, tol2 = 1e-04))

## Unit: milliseconds
##           expr
##  ADMMsigma(S = S, lam = 0.1, alpha = 1, tol1 = 1e-04, tol2 = 1e-04)
##           min       lq      mean    median      uq      max neval
##  40.48786 41.80211 45.23219 42.58309 44.18543 212.5204   100

# benchmark ADMMsigma - tolerance 1e-8
microbenchmark(ADMMsigma(S = S, lam = 0.1,
                        alpha = 1, tol1 = 1e-08, tol2 = 1e-08))

## Unit: milliseconds
##           expr
##  ADMMsigma(S = S, lam = 0.1, alpha = 1, tol1 = 1e-08, tol2 = 1e-08)
##           min       lq      mean    median      uq      max neval
##  185.7857 190.2324 198.1983 193.4001 198.3315 256.4238   100

# benchmark ADMMsigma CV - default
# parameter grid
microbenchmark(ADMMsigma(X), times = 5)

## Unit: seconds
##           expr      min       lq      mean    median      uq      max neval
##  ADMMsigma(X) 16.30647 16.30802 16.3881 16.40673 16.45475 16.46454    5

# benchmark ADMMsigma parallel CV
microbenchmark(ADMMsigma(X, cores = 3), times = 5)

## Unit: seconds
##           expr      min       lq      mean    median      uq
##  ADMMsigma(X, cores = 3) 12.95143 13.07855 13.29565 13.29788 13.54759
##           max neval
##  13.60278    5

# benchmark ADMMsigma CV - likelihood
# convergence criteria
microbenchmark(ADMMsigma(X, crit = "loglik"),

```

```

times = 5)

## Unit: seconds
##           expr      min      lq      mean  median
## ADMMsigma(X, crit = "loglik") 52.37604 52.58107 52.74499 52.68818
##           uq      max neval
## 52.77534 53.30433      5

```

3.4 Simulations

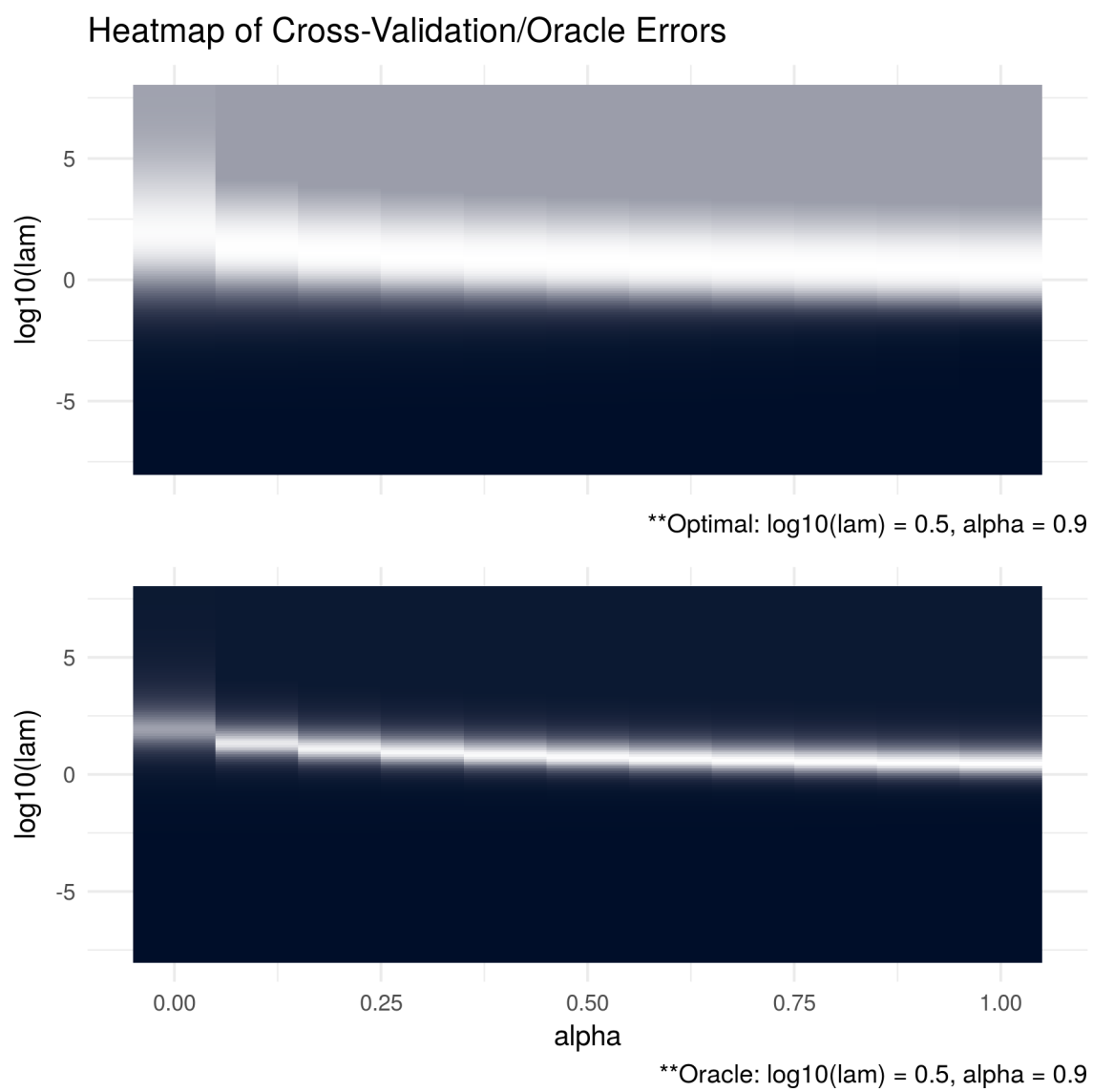


Figure 1:

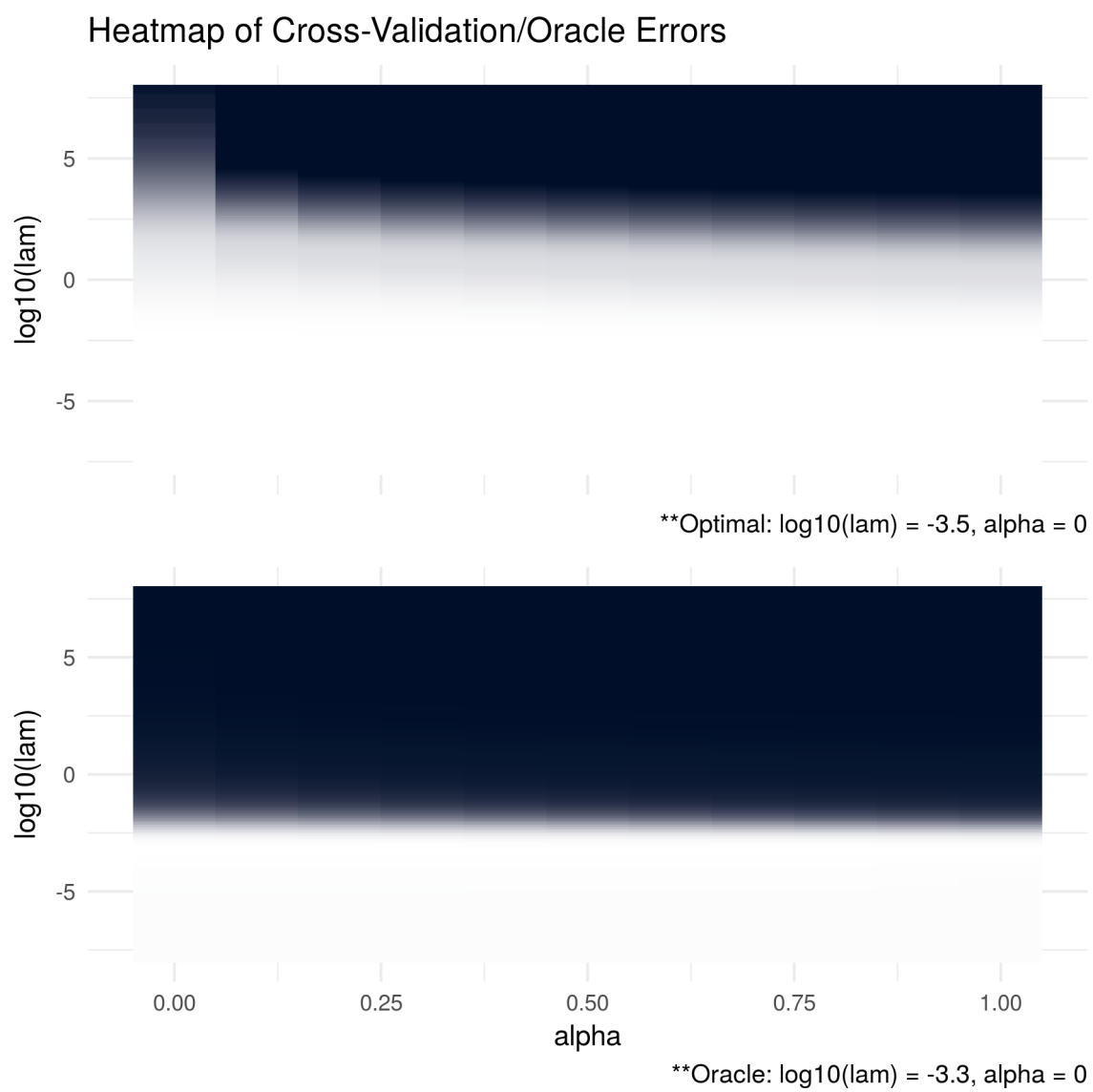


Figure 2:

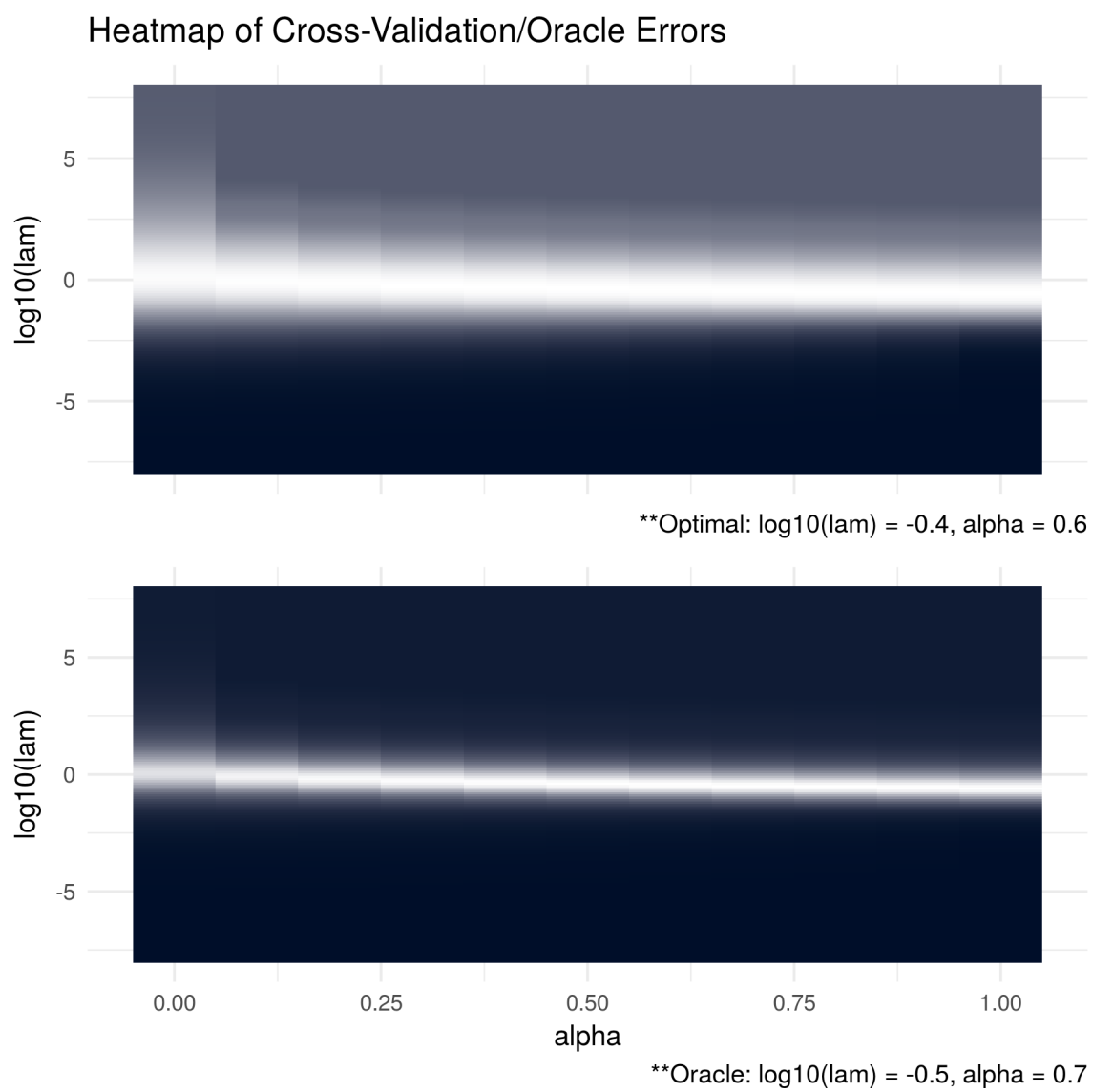


Figure 3:

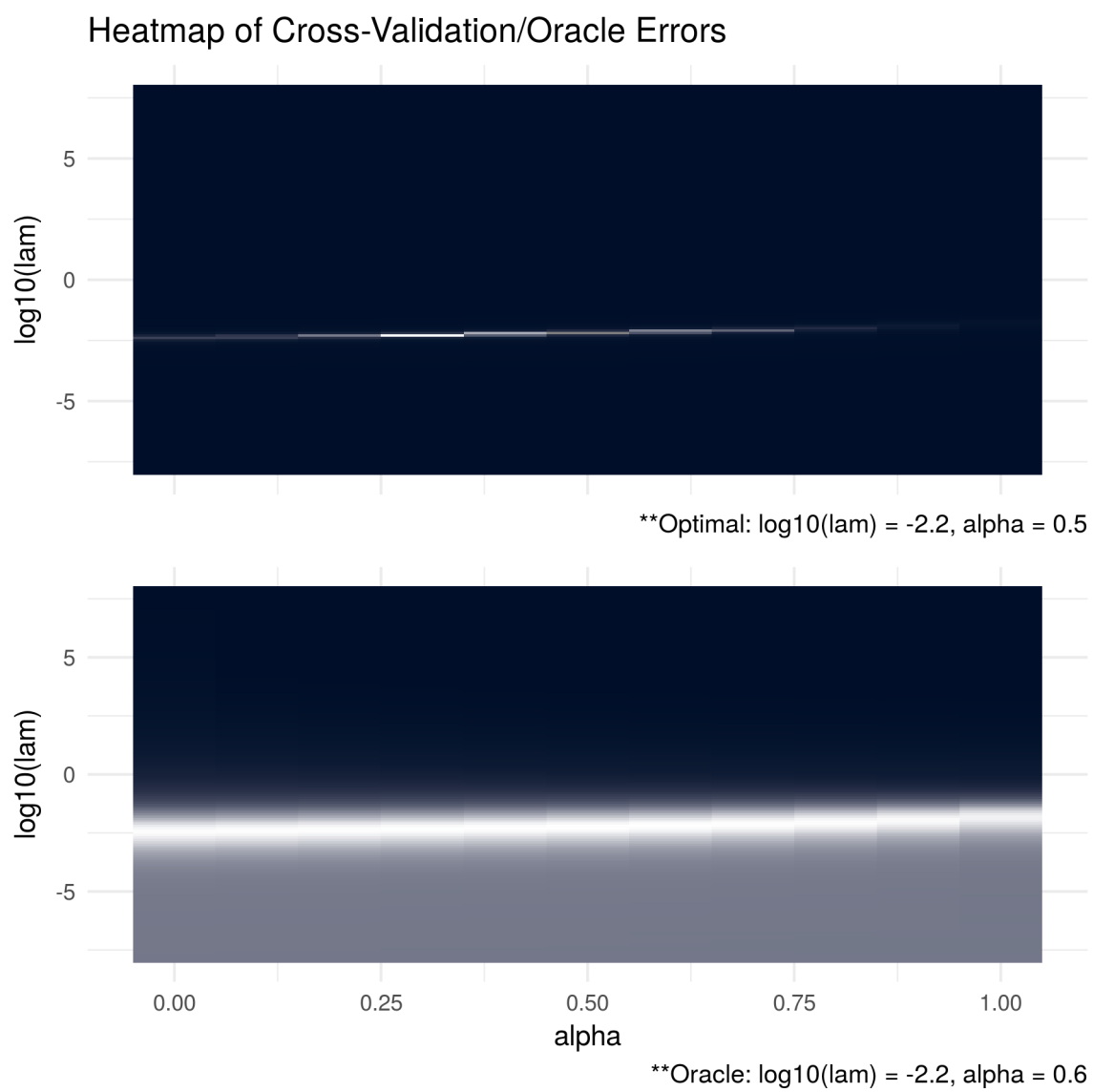


Figure 4:

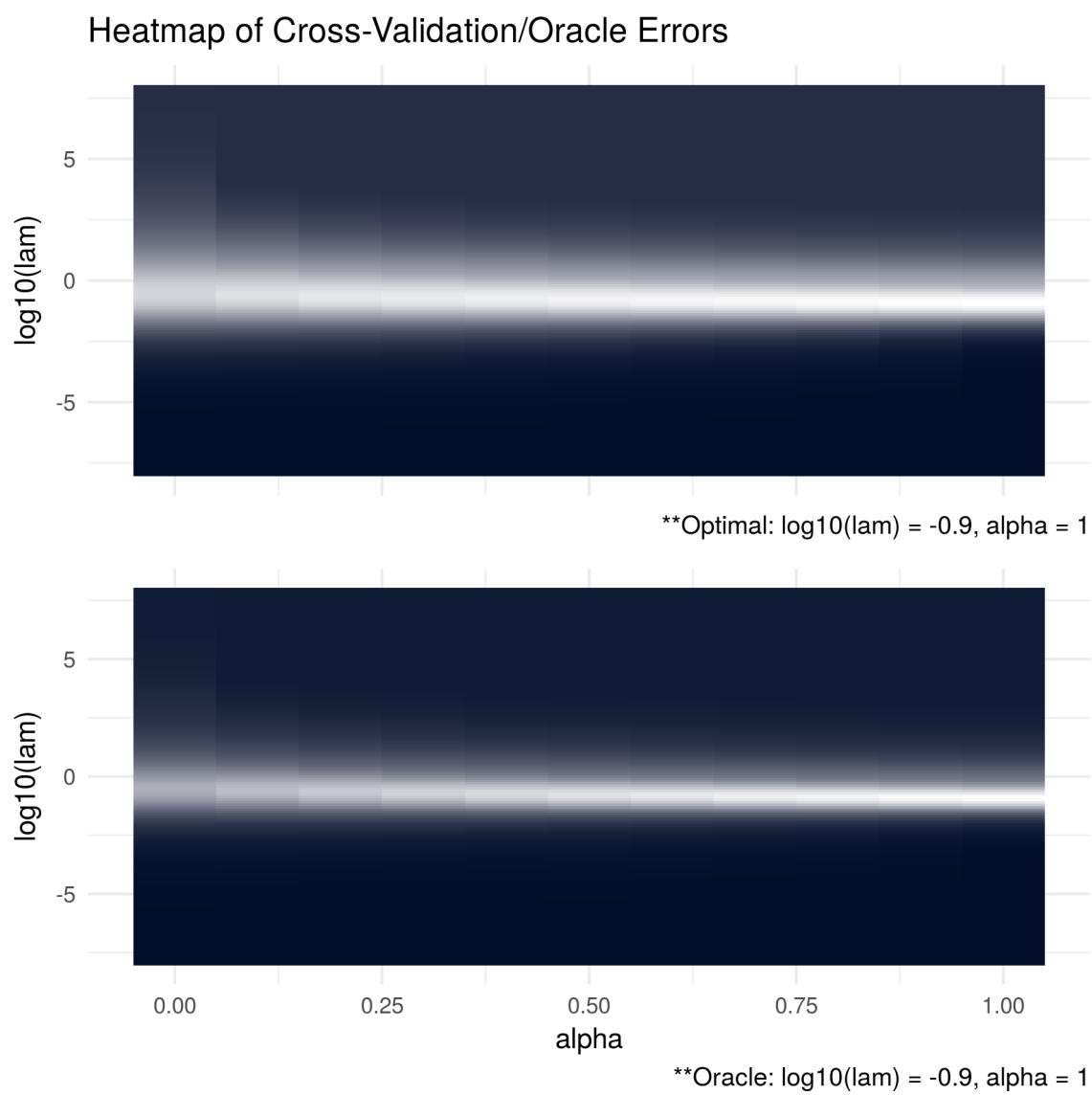


Figure 5:

References

- Boyd, Stephen, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, and others. 2011. “Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers.” *Foundations and Trends in Machine Learning* 3 (1). Now Publishers, Inc.: 1–122.
- Zou, Hui, and Trevor Hastie. 2005. “Regularization and Variable Selection via the Elastic Net.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 67 (2). Wiley Online Library: 301–20.