

---

# Regularized selection indices for breeding value prediction using hyper-spectral image data

Marco Lopez-Cruz<sup>1</sup>, Eric Olson<sup>1</sup>, Gabriel Rovere<sup>2,3,4</sup>, Jose Crossa<sup>6</sup>, Susanne Dreisigacker<sup>6</sup>,  
Suchismita Mondal<sup>6</sup>, Ravi Singh<sup>6</sup>, and Gustavo de los Campos<sup>3,4,5,\*</sup>

{<sup>1</sup>Department of Plant, Soil and Microbial Sciences, <sup>2</sup>Department of Animal Science,  
<sup>3</sup>Department of Epidemiology and Biostatistics, <sup>4</sup>Institute for Quantitative Health Science and  
Engineering, <sup>5</sup>Department of Statistics and Probability}, Michigan State University, USA  
<sup>6</sup>International Maize and Wheat Improvement Center (CIMMYT), Mexico

\*Corresponding author (e-mail: [gustavoc@msu.edu](mailto:gustavoc@msu.edu))

---

## Contents

<b>1</b>	<b>Selection index</b>	<b>2</b>
1.1	Canonical selection index . . . . .	2
<b>2</b>	<b>Regularized selection indices</b>	<b>2</b>
2.1	Principal components-based selection indices . . . . .	2
2.2	Penalized selection indices . . . . .	3
<b>3</b>	<b>Phenotypic and genetic parameters</b>	<b>3</b>
<b>4</b>	<b>Accuracy of indirect selection</b>	<b>4</b>
<b>5</b>	<b>Cross validation</b>	<b>4</b>
<b>6</b>	<b>Experimental data</b>	<b>4</b>
<b>7</b>	<b>Implementation</b>	<b>5</b>
7.1	Data downloading . . . . .	5
7.2	Heritability and variance components . . . . .	6
7.3	Training-testing partitions . . . . .	6
7.4	Genetic covariances estimation . . . . .	7
7.5	Single time-point selection indices . . . . .	7
7.5.1	Regression coefficients . . . . .	7
7.5.2	Accuracy of the index . . . . .	8
7.5.3	Effect of regularization on the accuracy of the index . . . . .	9
7.5.4	Optimal regularized vs canonical selection indices . . . . .	10
7.5.5	Comparison of the indices across time-points . . . . .	11
7.6	Multi time-point selection indices . . . . .	12
7.6.1	Regression coefficients . . . . .	12
7.6.2	Accuracy of the index . . . . .	13
7.6.3	Optimal multi vs single time-point selection indices . . . . .	14
7.7	Sparsity of the index . . . . .	15

# 1 Selection index

Recall the basic genetic model that decomposes phenotypic observations,  $\mathbf{y} = (y_1, \dots, y_n)'$ , as the sum of breeding values,  $\mathbf{u}_y = (u_{y_1}, \dots, u_{y_n})'$ , and environmental deviations,  $\mathbf{e}_y = (e_{y_1}, \dots, e_{y_n})'$ , as

$$y_i = u_{y_i} + e_{y_i} \quad (1)$$

Both  $\mathbf{u}_y$  and  $\mathbf{e}_y$  are assumed to have null means and  $\text{var}(\mathbf{u}_y) = \sigma_{u_y}^2 \mathbf{K}$ ,  $\text{var}(\mathbf{e}_y) = \sigma_{e_y}^2 \mathbf{I}$ , and  $\text{cov}(\mathbf{u}, \mathbf{e}') = \mathbf{0}$ , where  $\sigma_{u_y}^2$  and  $\sigma_{e_y}^2$  are the common genetic and residual variances, respectively;  $\mathbf{K}$  is the **kinship relationship matrix** and  $\mathbf{I}$  is an identity matrix.

A **selection index** (SI) (Hazel, 1943; Smith, 1936) is used to predict the genetic merit for a the selection goal ( $u_{y_i}$ , e.g., the genetic merit for grain yield of the  $i^{\text{th}}$  genotype) as a linear combination of  $p$  measured phenotypes,  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})'$ , of the form

$$\mathcal{I}_i = \mathbf{x}_i' \boldsymbol{\beta} \quad (2)$$

where  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$  is a vector of regression coefficients whose entries define the weights of each of the measured phenotypes in the SI.

## 1.1 Canonical selection index

In a **canonical SI** (CSI) the weights are derived by minimizing the expected squared deviation between the genetic merit for the selection goal and the index, that is

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \frac{1}{2} \mathbb{E} (u_{y_i} - \mathbf{x}_i' \boldsymbol{\beta})^2$$

Assuming that  $\mathbf{x}_i$  follows also model (1) and that their environmental effects ( $\mathbf{e}_{x_i}$ ) are orthogonal to  $u_{y_i}$ , the above problem is equivalent to

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left[ -\mathbf{G}_{x,y}' \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\beta}' \mathbf{P}_x \boldsymbol{\beta} \right] \quad (3)$$

where  $\mathbf{G}_{x,y} = \mathbb{E}(u_{y_i} \mathbf{x}_i)$  is a vector containing the genetic covariances between the selection target and each of the measured phenotypes, and  $\mathbf{P}_x = \mathbb{E}(\mathbf{x}_i \mathbf{x}_i')$  is the phenotypic (co)variance matrix of  $\mathbf{x}_i$ . The solution to this optimization problem is

$$\hat{\boldsymbol{\beta}} = \mathbf{P}_x^{-1} \mathbf{G}_{x,y} \quad (4)$$

The SI is by construction the **best linear predictor** (BLP) of the genetic merit for the selection goal; this property holds when  $\mathbf{G}_{x,y}$  and  $\mathbf{P}_x$  are known.

# 2 Regularized selection indices

## 2.1 Principal components-based selection indices

Principal components (PC) can be obtained from the singular value decomposition of the real-valued matrix,  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]'$  (the matrix containing all measured traits) which takes the form  $\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}'$ . Matrix  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]$  contains the left-singular vectors,  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_p]$  is the matrix with the right-singular vectors, and  $\mathbf{D} = \text{diag}(d_1, \dots, d_p)$  is a diagonal matrix with positive or zero elements. The PCs  $\mathbf{W} = \mathbf{X} \mathbf{V} = \mathbf{U} \mathbf{D}$  are then linear combinations of the measured phenotypes. A **principal components-based SI** (PC-SI) uses the first  $q$  PCs ( $\tilde{\mathbf{W}} = [\mathbf{w}_1, \dots, \mathbf{w}_q]$ ,  $q \leq p$ ) as 'measured phenotypes' in the SI, as

$$\mathcal{I}_i = \tilde{\mathbf{w}}'_i \boldsymbol{\gamma}^{(q)}$$

where  $\tilde{\mathbf{w}}'_i$  is a vector containing the scores for the  $i^{th}$  observation on the first  $q$  PCs and  $\boldsymbol{\gamma}^{(q)}$  are their corresponding regression coefficients. The solution to the optimization problem takes the form  $\hat{\boldsymbol{\gamma}}^{(q)} = (n-1)(\tilde{\mathbf{D}}^2)^{-1} \mathbf{G}_{\tilde{\mathbf{w}},y}$ , where  $\tilde{\mathbf{D}}$  contains only the first  $q$  elements of  $\mathbf{D}$ , and  $\mathbf{G}_{\tilde{\mathbf{w}},y}$  is a vector containing the genetic covariances between each of the top  $q$  PCs and the selection objective. This solution can be mapped to coefficients for the measured traits, as in expression (2), using

$$\hat{\boldsymbol{\beta}}^{(q)} = (n-1) \tilde{\mathbf{V}} (\tilde{\mathbf{D}}^2)^{-1} \mathbf{G}_{\tilde{\mathbf{w}},y} \quad (5)$$

where  $\tilde{\mathbf{V}}$  is the matrix containing only the first  $q$  right-singular vectors.

## 2.2 Penalized selection indices

In a penalized regression, regularization is achieved by including in the objective function a penalty on model complexity. A **penalized selection index (PSI)** is obtained by considering an optimization problem as in (3), as follows:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left[ -\mathbf{G}'_{x,y} \boldsymbol{\beta} + \frac{1}{2} \boldsymbol{\beta}' \mathbf{P}_x \boldsymbol{\beta} + \lambda \cdot J(\boldsymbol{\beta}) \right] \quad (6)$$

where  $\lambda$  is a penalty parameter and  $J(\boldsymbol{\beta}_i)$  is a penalty function. Commonly used penalties include the L2 ( $\|\boldsymbol{\beta}_i\|_2^2 = \sum_{j=1}^n \beta_{ij}^2$ ) and L1 ( $\|\boldsymbol{\beta}_i\|_1 = \sum_{j=1}^n |\beta_{ij}|$ ) norms (Fu, 1998) either alone or in a combination of both. We considered a penalty function as

$$J(\boldsymbol{\beta}) = \frac{1}{2} (1 - \alpha) \sum_{j=1}^n \beta_j^2 + \alpha \sum_{j=1}^n |\beta_j|$$

where  $\alpha$  is a weighting factor. This penalty is used in an **elastic-net-type** regression (Zou and Hastie, 2005) that combines the shrinkage-inducing feature of a **ridge-regression** (Hoerl and Kennard, 1970), that uses the L2-norm alone, and the variable selection feature of a **LASSO** regression (Tibshirani, 1996), that uses the L1-norm alone. With no penalization ( $\lambda = 0$ ), the solution for (6) is equivalent to that of the canonical SI in (4).

The ridge-regression (**L2-PSI**) and LASSO (**L1-PSI**) types are special cases of the elastic-net-type PSI (**EN-PSI**) when  $\alpha = 0$  and  $\alpha = 1$ , respectively. A closed-form solution for the regression coefficients can be found only for the L2-PSI (Hastie et al., 2009). In this case, the solution is

$$\hat{\boldsymbol{\beta}}^{L2} = (\mathbf{P}_x + \lambda \mathbf{I})^{-1} \mathbf{G}_{x,y}$$

where  $\mathbf{I}$  is a  $p \times p$  identity matrix. For the EN-PSI ( $0 < \alpha \leq 1$ ) solutions are obtained using iterative algorithms such as **least angle regression** (LARS, Efron et al., 2004) or **coordinate descent** (Friedman et al., 2007) for different values of the parameters  $\alpha$  and  $\lambda$ . These combinations of the values of  $\alpha$  and  $\lambda$  will result in different EN-PSI from which an optimal index can be obtained such as the prediction accuracy is maximum.

## 3 Phenotypic and genetic parameters

The population phenotypic (co)variance matrix  $\mathbf{P}_x$  among the measured phenotypes is estimated using the unbiased sample (co)variance matrix given by

$$\hat{\mathbf{P}}_x = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})'$$

where  $\bar{\mathbf{x}}$  is the vector containing the sample mean of the measured phenotypes. For centered data, this reduces to  $\hat{\mathbf{P}}_x = \frac{1}{n-1} \mathbf{X}'\mathbf{X}$ , where  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]'$  is the matrix containing all measured traits.

The genetic covariance ( $\mathbf{G}_{y,x_j}$ ) between the target trait and the  $j^{\text{th}}$  measured trait ( $j = 1, \dots, p$ ) is estimated using a sequence of uni-variate genetic models as in equation (1). First, the model is fitted for the target trait as response, then for each of the measured traits, and then for the sum of the target trait and each of the measured traits. The genetic covariances between the target trait and the measured traits are then estimated using

$$\hat{\mathbf{G}}_{y,x_j} = \frac{1}{2} \left( \hat{\sigma}_{u_{y+x_j}}^2 - \hat{\sigma}_{u_y}^2 - \hat{\sigma}_{u_{x_j}}^2 \right) \quad (7)$$

where  $\hat{\sigma}_{u_y}^2$ ,  $\hat{\sigma}_{u_{x_j}}^2$  and  $\hat{\sigma}_{u_{y+x_j}}^2$  are the estimated genetic variances for the target trait, the  $j^{\text{th}}$  measured trait, and the sum of the target trait and the  $j^{\text{th}}$  measured trait, respectively.

## 4 Accuracy of indirect selection

The accuracy of the index to be used for indirect selection,  $\text{Acc}(\mathcal{I})$ , is defined as the correlation between the index ( $\mathcal{I}_i$ ) and the breeding values ( $u_i$ ). This parameter is equal to the product of the square root of the heritability of the SI,  $h_{\mathcal{I}}$ , times the genetic correlation between the SI and the selection target,  $\text{cor}(u_{\mathcal{I}_i}, u_{y_i})$ , this is

$$\text{Acc}(\mathcal{I}) = \text{cor}(\mathcal{I}_i, u_i) = \text{cor}(u_{\mathcal{I}_i}, u_{y_i}) h_{\mathcal{I}}$$

The efficiency of indirect selection relative to mass phenotypic selection ( $RE$ ) is calculated as the ratio between the accuracy of indirect selection and the accuracy of direct selection,  $h_y$ , this is

$$RE = \frac{h_{\mathcal{I}}}{h_y} \text{cor}(u_{\mathcal{I}_i}, u_{y_i})$$

where  $h_y$  is the square root of the heritability of the selection goal ( $y$ ), calculated as

$$h_y^2 = \frac{\sigma_{u_y}^2}{\sigma_{u_y}^2 + \sigma_{e_y}^2}$$

## 5 Cross validation

To avoid bias in the estimation, the accuracy of indirect selection must be estimated using data that was not used to derive the coefficients of the index. To this end, a **training-testing partition** procedure is implemented as follows: (i) data is partitioned into training (TRN) and testing (TST) sets, (ii) the coefficients of the SI,  $\beta$ , are derived in the training set, (iii) these coefficients are applied to data of the testing set to derive an index (using (2)), and (iv) parameters  $\text{cor}(u_{\mathcal{I}_i}, u_{y_i})$ ,  $h_{\mathcal{I}}$ , and  $\text{Acc}(\mathcal{I})$  are estimated in the testing set.

## 6 Experimental data

The data set consists of 1,092 inbred wheat lines grouped into 39 trials and grown during the 2013- 2014 season at the Norman Borlaug experimental research station in Ciudad Obregon, Sonora, Mexico. Each trial consisted of 28 breeding lines that were arranged in an alpha-lattice design with three replicates and six sub-blocks. The trials were grown in four different environments:

- E1: *Flat-Drought*. (1092 genotypes, 3274 records) Sowing in flat with irrigation of 180 mm through drip system
- E2: *Bed-2IR*. (1090 genotypes, 3267 records) Sowing in bed with 2 irrigations approximately 250 mm)
- E3: *Bed-5IR*. (1092 genotypes, 3274 records) Bed sowing with 5 normal irrigations
- E4: *Bed-EHeat*. (1013 genotypes, 2858 records) Bed sowing 30 days before optimal planting date with 5 normal irrigations of approximately 500 mm)

Records of grain yield (YLD) were collected as the total plot yield after maturity. Measurements for days to heading (DTH), days to maturity (DTM), and plant height (PH) were recorded only in the first replicate at each trial. Reflectance phenotypic data were collected from the fields using both infrared (A600 series Infrared camera, FLIR, Wilsonville, OR) and hyper-spectral (A-series, Micro-Hyperspec, VNIR Headwall Photonics, Fitchburg, MA) cameras mounted on a Piper PA-16 Clipper aircraft on 9 different dates (time-points) between January 10<sup>th</sup> and March 27<sup>th</sup>, 2014. During each flight, data from 250 wavelengths ranging from 392 to 850 nm were collected for each pixel in the pictures. The average reflectance of all the pixels for each wavelength was calculated from each of the geo-referenced trial plots and reported as each line reflectance.

Grain yield and image data were pre-adjusted using mixed-effects model that accounted for genotype, trial, replicate, and sub-block. No pre-adjusting was made for DTH, DTM, and PH traits. Pre-adjusted phenotypes were obtained by subtracting from the phenotypic record the mean plus BLUPs of trial, replicate, and sub-block effects.

## 7 Implementation

The above described data will be used to generate selection indices for grain yield (as selection target) using all the wavelengths collected from cameras (measured phenotypes). Analyses will be implemented in R software using the SFSI R-package.

### 7.1 Data downloading

Both phenotypic and reflectance data can be downloaded from the `github` version of the SFSI R-package at [https://github.com/MarcoLopez/SFSI\\_data/tree/master/data](https://github.com/MarcoLopez/SFSI_data/tree/master/data). Each `.RData` file corresponds to one environment and contains objects `Y`, `WL`, and `VI`. Matrix `Y` contains pre-adjusted observations for YLD, DTH, DTM, and PH. Object `WL` is a list type containing the hyper-spectral image data for all 9 time-points. Element `WL[[j]]`,  $j = 1, \dots, 9$ , is a matrix with 250 columns (wavelengths) and rows matched with those of matrix `Y`. Likewise, the object `VI` is a list with 9 elements being two-columns matrices containing green and red NDVI. **Box 1** below shows how to prepare data for a single environment which will be used in later analyses.

#### Box 1. Data downloading

```
rm(list = ls())
setwd("/Users/epidemiology/Dropbox/MSU/PBGB/projects/PSI/pipeline")
site <- "https://github.com/MarcoLopez/SFSI_data/blob/master/data"
filename <- "wheat.E1.RData"

# Download files
dir.create("data", recursive=TRUE)
url <- paste0(site,"/",filename,"?raw=true")
download.file(url, destfile= paste0("./data/",filename), mode = "wb")
```

## 7.2 Heritability and variance components

Heritability can be calculated from variance components obtained by fitting the model (1) for the whole data. Code below illustrates how to fit this model using the `solveMixed` function from the `SFSI` R-package.

### Box 2. Variance components calculation

```
library(SFSI)

# Load data
load(paste0("data/",filename))
Y$gid <- factor(as.character(Y$gid))
y <- as.vector(scale(Y[, "YLD"]))

# Fit model
Z <- model.matrix(~0+gid,data=Y)
evdZZt <- eigen(tcrossprod(Z)) # Decomposition of ZZ'
fm0 <- solveMixed(y,U=evdZZt$vectors,d=evdZZt$values,BLUP=FALSE)
c(fm0$varU,fm0$varE,fm0$h2)

dir.create("output", recursive=TRUE)
save(fm0,evdZZt,file="output/varComps.RData")
```

## 7.3 Training-testing partitions

Code in **Box 3** below illustrates how to split data into TRN and TST sets. In this example,  $\frac{2}{3}$  of the 29 trials (26 trials,  $n_{TRN} \approx 2,184$  observations) will be randomly assigned to the training set and the remaining  $\frac{1}{3}$  (13 trials,  $n_{TST} \approx 1,092$ ) to the testing set.

The output will be a 'list' type of length `nPart` containing indices indicating which observations are assigned to testing sets. The object will be saved in the file `partitions.RData` and will be used later for single time-point and across time-points analyses.

For demonstration purposes, only `nPart=5` will be run throughout all the analyses.

### Box 3. Create testing set partitions

```
trials <- as.numeric(as.character(Y$trial))
nTrial <- length(unique(trials))
nTST <- ceiling((1/3)*nTrial) # Number of trials in TST set

nPart <- 5 # Number of partitions
seeds <- round(seq(1E3, .Machine$integer.max, length = nPart))

partitions <- vector("list",length=nPart) # Object to store partitions

for(k in 1:length(partitions))
{
  set.seed(seeds[k])
  tst <- sample(1:nTrial,nTST,replace=FALSE)
  partitions[[k]] <- which(trials %in% tst)
}
save(partitions,file="output/partitions.RData")
```

## 7.4 Genetic covariances estimation

The genetic covariances between the grain yield and each of the wavelengths will be calculated using a sequence of models as described in (7) implemented in the `getGenCov` function from the **SFSI** package. Code in **Box 4** below shows how to calculate genetic and phenotypic covariances within time-point. The time-point can be specified through variable `tp` thus it can be changed to calculate covariances for all the 9 time-points. Calculations are performed using data from training set only for each partition. Results are stored in matrices `gencov` and `phencov`, and saved in the file `covariances.RData`.

### Box 4. Genetic covariances within time-point

```
tp <- 9      # Time-point
load("output/partitions.RData")

gencov <- phencov <- c()  # Matrices to store covariances

for(k in 1:length(partitions))
{
  cat("  partition = ",k,"of",length(partitions),"\\n")
  indexTST <- partitions[[k]]
  indexTRN <- (1:nrow(Y))[-indexTST]

  # Training set
  xTRN <- scale(WL[[tp]][indexTRN,])
  yTRN <- as.vector(scale(Y[indexTRN,"YLD"]))

  genotype <- factor(as.character(Y[indexTRN,"gid"]))
  Z <- model.matrix(~ 0 + genotype)  # Design matrix for genotypes
  ZZt <- tcrossprod(Z)

  # Get genetic variances and covariances
  fm <- getGenCov(y1=yTRN,y2=xTRN,K=ZZt,mc.cores=5,scale=FALSE)

  gencov <- cbind(gencov,fm$covU)
  phencov <- cbind(phencov,fm$covU + fm$covE)
}
dir.create(paste0("output/timepoint_",tp), recursive=TRUE)
save(gencov,phencov,file=paste0("output/timepoint_",tp,"/covariances.RData"))
```

## 7.5 Single time-point selection indices

Canonical, L1-penalized and PC-based SIs will be fitted using data from a single time-point. The regression coefficients are to be obtained from training data within each partition previously created.

### 7.5.1 Regression coefficients

Code in **Box 5a** can be used to estimate regression coefficients using training data within each partition following equations (4) (for the CSI), and (5) (with different number of PCs for the PC-SI). Coefficients for L1-PSI are estimated as in expression (6) whose solutions are found using the `lars2` function from the **SFSI** R-package. This function implements the LARS algorithm that gives solutions for an entire path of values of parameter  $\lambda$ .

Time-point can be specified through variable `tp` thus it can be changed to get results for all time-points (1, 2, ..., 9). Results are stored in matrices `betaCSI`, `betaPCSI` and `betaL1PSI`, and saved in the file `coefficients.RData`. Calculations are performed across all partitions previously created.

**Box 5a.** Regression coefficients within time-point

```

tp <- 9      # Time-point
load(paste0("output/timepoint_",tp,"/covariances.RData"))
load("output/partitions.RData")

# Objects to store regression coefficients
betaCSI <- betaPCSI <- betaL1PSI <- vector("list",length(partitions))

for(k in 1:length(partitions))
{
  cat("  partition = ",k,"of",length(partitions),"\\n")
  indexTST <- partitions[[k]]
  indexTRN <- (1:nrow(Y))[-indexTST]

  # Training set
  xTRN <- scale(WL[[tp]][indexTRN,])
  VARx <- var(xTRN)
  EVDx <- eigen(VARx)

  # -- Canonical SI --
  VARinv <- EVDx$vectors %*% diag(1/EVDx$values) %*% t(EVDx$vectors)
  betaCSI[[k]] <- VARinv %*% gencov[,k]

  # -- PC-SI --
  VARinv <- diag(1/EVDx$values)
  gamma <- as.vector(VARinv %*% t(EVDx$vectors) %*% gencov[,k])
  beta <- apply(EVDx$vectors %*% diag(gamma),1,cumsum)
  betaPCSI[[k]] <- data.frame(I(beta),df=1:nrow(beta),lambda=NA)

  # -- L1-PSI --
  fm <- lars2(VARx, gencov[,k])
  beta <- as.matrix(fm$beta)[-1,]
  betaL1PSI[[k]] <- data.frame(I(beta),df=fm$df[-1],lambda=fm$lambda[-1])
}
save(betaCSI,betaPCSI,betaL1PSI,file=paste0("output/timepoint_",tp,"/coefficients.RData"))

```

**7.5.2 Accuracy of the index**

The regression coefficients above calculated can be applied to data from testing set to derive selection indices (as in equation (2)). Code in **Box 5b** can be used to (i) compute the SIs for the testing data, and then (ii) calculate the accuracy of the index, the heritability of the index, and genetic correlation of the index with grain yield. The genetic covariances are obtained using again the approach in (7) implemented in the `getGenCov` function.

Results for the three indices (CSI, PC-SI and L1-PSI) are stored in the object matrix `accSI` and saved in the file `accuracies.RData`. Calculations are performed across all partitions previously created.



**Box 5b.** Accuracy of selection within time-point

```

tp <- 9      # Time-point
load(paste0("output/timepoint_",tp,"/coefficients.RData"))
load("output/partitions.RData")

accSI <- c()  # Object to store accuracy components

for(k in 1:length(partitions))
{
  cat("  partition = ",k,"of",length(partitions),"\\n")
  indexTST <- partitions[[k]]

  # Testing set
  xTST <- scale(WL[[tp]][indexTST,])
  yTST <- as.vector(scale(Y[indexTST,"YLD"]))
  genotype <- factor(as.character(Y[indexTST,"gid"]))
  Z <- model.matrix(~ 0 + genotype)    # Design matrix for genotypes

  # Calculate the indices
  CSI <- xTST %*% betaCSI[[k]]
  PCSI <- tcrossprod(xTST, betaPCSI[[k]]$beta)
  L1PSI <- tcrossprod(xTST, betaL1PSI[[k]]$beta)

  # Fit genetic models
  ZZt <- tcrossprod(Z)
  SI <- as.matrix(data.frame(I(CSI),I(PCSI),I(L1PSI)))
  fm <- getGenCov(y1=yTST,y2=SI,K=ZZt)

  # Retrieve accuracy components
  h <- sqrt(fm$varU2/(fm$varU2 + fm$varE2))
  gencor <- fm$covU/sqrt(fm$varU1*fm$varU2)
  accuracy <- gencor*h

  df <- c(max(betaL1PSI[[k]]$df),betaPCSI[[k]]$df,betaL1PSI[[k]]$df)
  lambda <- c(min(betaL1PSI[[k]]$lambda),betaPCSI[[k]]$lambda,betaL1PSI[[k]]$lambda)
  accSI <- rbind(accSI,data.frame(rep=k,SI=colnames(SI),h,gencor,accuracy,df,lambda))
}
save(accSI,file=paste0("output/timepoint_",tp,"/accuracies.RData"))

```

**7.5.3 Effect of regularization on the accuracy of the index**

The regularized selection indices were calculated for  $q = 1, 2, \dots, 250$  PCs for the PC-SI and for decreasing values,  $\lambda_1, \dots, \lambda_{250}$ , of the path of the penalization for the L1-PSI. In the later case, decreasing  $\lambda$  will give increasing number of active number of predictors in the index. Code in **Box 6** below can be used to create a plot depicting the evolution of accuracy (and its components) over values of regularization (number of PCs or number of active predictors). Results for the CSI corresponds to the case of a PC-SI with 250 PCs and a L1-PSI with 250 active predictors. These cases are shown at the right-most results in the plots.

**Box 6.** Effect of regularization

```

library(ggplot2); library(reshape); library(ggpubr)

tp <- 9      # Time-point
load(paste0("output/timepoint_",tp,"/accuracies.RData"))

accSI <- split(accSI,as.character(accSI$SI))
accSI <- data.frame(do.call(rbind,lapply(accSI,function(x)apply(x[, -c(1,2)],2,mean))))
accSI$SI <- unlist(lapply(strsplit(rownames(accSI),"\\"),function(x)x[1]))

# Plot of PC-SI
dat <- melt(accSI[accSI$SI=="PCSI",],id=c("df","lambda","SI"))
plot1 <- ggplot(dat[dat$df>1,],aes(df,value,group=variable,color=variable)) + theme_bw() +
  geom_line(size=1) + labs(title="PC-SI",y="correlation",x="Number of PCs")

# Plot of L1-PSI
dat <- melt(accSI[accSI$SI=="L1PSI",],id=c("df","lambda","SI"))
plot2 <- ggplot(dat[dat$df>1,],aes(df,value,group=variable,color=variable)) + theme_bw() +
  geom_line(size=1) + labs(title="L1-PSI",y="correlation",x="Number of active predictors")

ggarrange(plot1,plot2,nrow=1,common.legend=TRUE,legend="right")

```

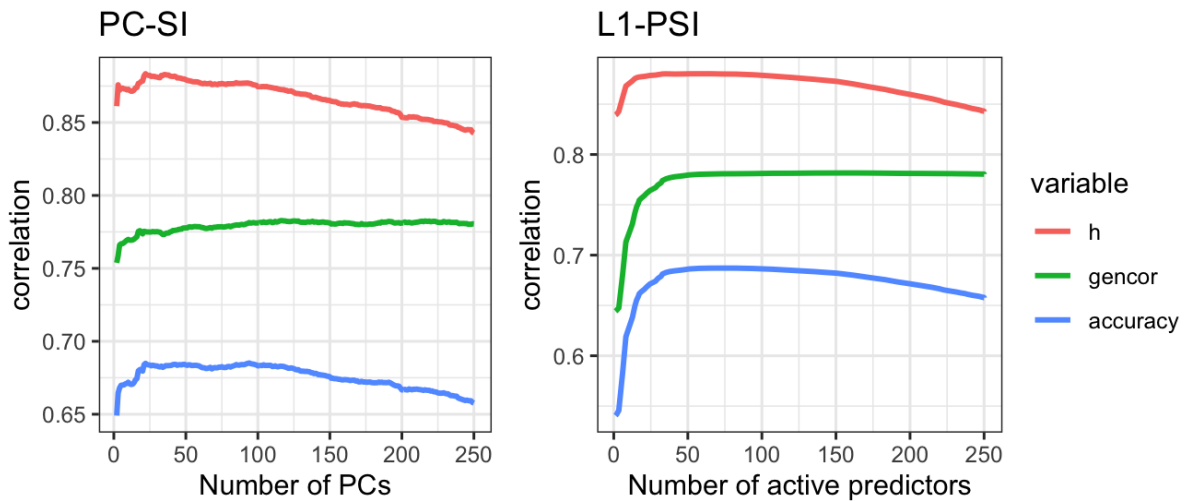


Figure 1: Accuracy of selection (average across 10 TRN-TST partitions) and its components of the index versus the number of principal components entering in the PC-SI (left) and the number of active predictors in the L1-PSI (right). Timepoint 9. Environment E1

#### 7.5.4 Optimal regularized vs canonical selection indices

An optimal regularized selection index can be chosen by selecting an optimal regularization (number of PCs or value of  $\lambda$ ) such that the accuracy of selection of the resulting index in the testing set is maximum. Code in **Box 7** below can be used to select the optimal PC-SI and L1-PSI, and to compare the accuracy of selection with that of the CSI (non-regularized). Results are reported as an average across all partitions along with a 95% confidence interval.

**Box 7. Optimal selection index**

```

tp <- 9      # Time-point
load(paste0("output/timepoint_",tp,"/accuracies.RData"))

accSI$SI <- unlist(lapply(strsplit(as.character(accSI$SI),"\\"),function(x)x[1]))
accSI <- split(accSI,paste(accSI$rep,"_",accSI$SI))
accSI <- do.call(rbind,lapply(accSI,function(x)x[which.max(x$accuracy),]))

dat = aggregate(accuracy ~ SI, mean, data=accSI)
dat$sd = aggregate(accuracy ~ SI,sd,data=accSI)$accuracy
dat$n = aggregate(accuracy ~ SI,length,data=accSI)$accuracy
dat$se = qnorm(0.975)*dat$sd/sqrt(dat$n)

ggplot(dat,aes(SI,accuracy,ymin=accuracy-se,ymax=accuracy+se)) + theme_bw() +
  geom_bar(stat="identity",width=0.5,fill="orange") + geom_errorbar(width=0.2) +
  geom_text(aes(label=sprintf("%.3f",accuracy)),y=min(dat$accuracy)*0.75)+labs(x="")

```

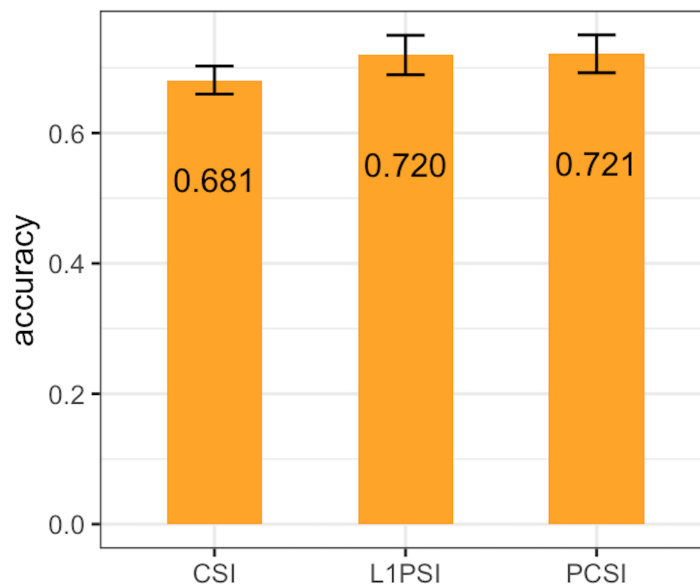


Figure 2: Accuracy of selection (average across 10 TRN-TST partitions) of the CSI, and the optimal PC-SI and L1-PSI. Timepoint 9. Environment E1

### 7.5.5 Comparison of the indices across time-points

The accuracy of indirect selection for the optimal (i.e., the one with the highest accuracy of indirect selection) L1-PSI and PC-SI can be assessed across time-points and compared with a canonical SI. Code in **Box 8** below can be used to calculate the optimal PC-SI and L1-PSI, and to compare their accuracy of selection with that of the CSI across all time-points.

To this end, all the within time-point analyses in **Box 4b**, **Box 5a** and **Box 5b** must be previously performed for all the 9 time-points. This can be achieved using, for instance, a loop as: `for(tp in 1:9){ }`.

**Box 8. Accuracy across time-points**

```

AccSI <- c()
for(tp in 1:9)
{ load(paste0("output/timepoint_",tp,"/accuracies.RData"))
  AccSI <- rbind(AccSI,data.frame(Timepoint=tp,accSI))
}

AccSI$SI <- unlist(lapply(strsplit(as.character(AccSI$SI),"\\"),function(x)x[1]))
AccSI <- split(AccSI,paste(AccSI$Timepoint,"_",AccSI$rep,"_",AccSI$SI))
AccSI <- do.call(rbind,lapply(AccSI,function(x)x[which.max(x$accuracy),]))
AccSI$Timepoint <- factor(as.character(AccSI$Timepoint))

dat = aggregate(accuracy ~ SI + Timepoint,mean,data=AccSI)
dat$sd = aggregate(accuracy ~ SI + Timepoint,sd,data=AccSI)$accuracy
dat$n = aggregate(accuracy ~ SI + Timepoint,length,data=AccSI)$accuracy
dat$se = qnorm(0.975)*dat$sd/sqrt(dat$n)

ggplot(dat,aes(Timepoint,accuracy,color=SI,group=SI,ymin=accuracy-se,ymax=accuracy+se))+
  geom_line() + geom_point() + theme_bw() + geom_errorbar(width=0.2)

```

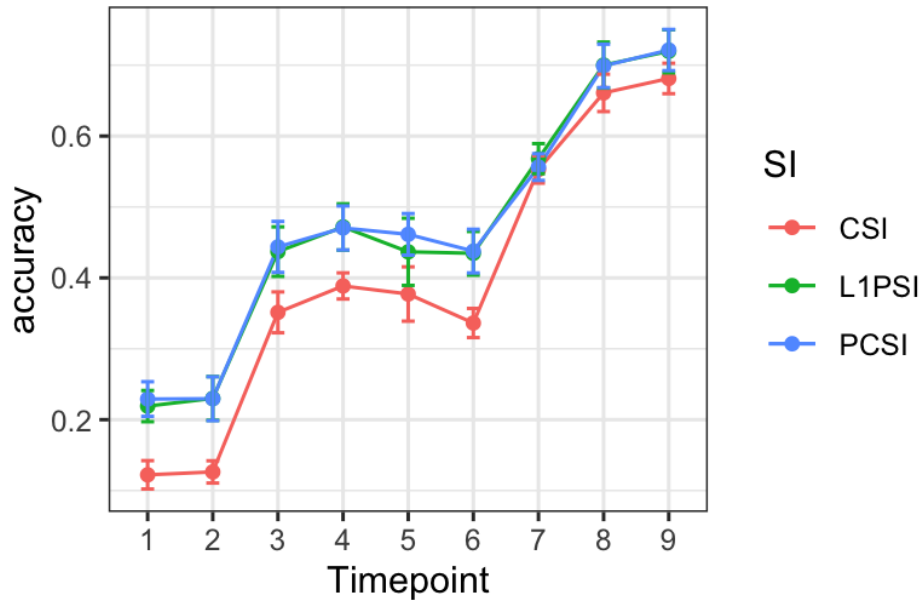


Figure 3: Accuracy of selection (average across 10 TRN-TST partitions) of the CSI, and the optimal PC-SI and L1-PSI. Timepoint 9. Environment E1

## 7.6 Multi time-point selection indices

Selection indices can be calculated to include reflectance data from multiple time-points. In this case, the vector of measured phenotypes  $\mathbf{x}_i$  will contain 2,250 traits corresponding to 250 wavebands measured at each of 9 time-points.

### 7.6.1 Regression coefficients

Code in **Box 9a** shows how to calculate coefficients for the PC-SI and L1-PSI using data from all 9 time-points. To this end, the genetic covariances of all 2,250 traits with grain yield are needed, thus, all the within time-point genetic covariances must be previously calculated for all the 9 time-points using code in **Box 4b**. Results are saved in file `coefficients.RData` in folder `output/multi_timepoint`.

**Box 9a.** Regression coefficients multi time-points

```

Gencov <- c()
for(tp in 1:9)
{ load(paste0("output/timepoint_",tp,"/covariances.RData"))
  Gencov <- rbind(Gencov,gencov)
}
load("output/partitions.RData")

# Objects to store regression coefficients
betaPCSI <- betaL1PSI <- vector("list",length(partitions))

for(k in 1:length(partitions))
{ cat(" partition = ",k,"of",length(partitions),"\\n")
  indexTST <- partitions[[k]]
  indexTRN <- (1:nrow(Y))[-indexTST]

  # Training set
  xTRN <- scale(do.call(cbind,WL)[indexTRN,])
  VARx <- var(xTRN)
  EVDx <- eigen(VARx)

  # -- PC-SI --
  gamma <- t(sweep(EVDx$vectors,2,1/EVDx$values,FUN="*")) %*% Gencov[,k]
  beta <- apply(sweep(EVDx$vectors,2,as.vector(gamma),FUN="*"),1,cumsum)[1:500,]
  betaPCSI[[k]] <- data.frame(I(beta),df=1:nrow(beta),lambda=NA)

  # -- L1-PSI --
  fm <- solveEN(VARx, Gencov[,k],nLambda=100,maxIter=200,tol=5E-4)
  beta <- as.matrix(fm$beta)[-1,]
  betaL1PSI[[k]] <- data.frame(I(beta),df=fm$df[-1],lambda=fm$lambda[-1])
}
dir.create("output/multi_timepoint", recursive=TRUE)
save(betaPCSI,betaL1PSI,file="output/multi_timepoint/coefficients.RData")

```

In the above code, the `solveEN` function included in the `SFSI` R-package was used for convenience to calculate the regression coefficients. `solveEN` function will yield solutions for a given grid of values of  $\lambda$  (covering all the space of possible values) rather than all the 2,250 solutions for the entire  $\lambda$  path provided by `lars2` function.

## 7.6.2 Accuracy of the index

As for the single time-point, the regression coefficients are applied to data from testing set to derive selection indices. Code in **Box 9b** can be used to (i) compute the multi time-point PC-SI and L1-PSI for the testing data, and then (ii) calculate the accuracy of the index, the heritability of the index, and genetic correlation of the index with grain yield. The genetic covariances are obtained using the `getGenCov` function.

Results for the two indices (PC-SI and L1-PSI) are stored in the object matrix `AccSI` and saved in the file `accuracies.RData` in folder `output/multi_timepoint`. Calculations are performed across all partitions previously created.

**Box 9b.** Accuracy of selection multi time-points

```

load("output/partitions.RData")
load("output/multi_timepoint/coefficients.RData")

AccSI <- c()      # Objects to store accuracy components

for(k in 1:length(partitions))
{  cat("  partition = ",k,"of",length(partitions),"\\n")
   indexTST <- partitions[[k]]

   # Testing set
   xTST <- scale(do.call(cbind,WL)[indexTST,])
   yTST <- as.vector(scale(Y[indexTST,"YLD"]))
   genotype <- factor(as.character(Y[indexTST,"gid"]))
   Z <- model.matrix(~ 0 + genotype)

   # Calculate the indices
   PCSI_multi <- tcrossprod(xTST, betaPCSI[[k]]$beta)
   L1PSI_multi <- tcrossprod(xTST, betaL1PSI[[k]]$beta)

   # Fit genetic models
   ZZt <- tcrossprod(Z)
   SI <- as.matrix(data.frame(I(PCSI_multi),I(L1PSI_multi)))
   fm <- getGenCov(y1=yTST,y2=SI,K=ZZt)

   # Retrieve accuracy components
   h <- sqrt(fm$varU2/(fm$varU2 + fm$varE2))
   gencor <- fm$covU/sqrt(fm$varU1*fm$varU2)
   accuracy <- gencor*h

   df <- c(betaPCSI[[k]]$df,betaL1PSI[[k]]$df)
   lambda <- c(betaPCSI[[k]]$lambda,betaL1PSI[[k]]$lambda)
   AccSI <- rbind(AccSI,data.frame(rep=k,SI=colnames(SI),h,gencor,accuracy,df,lambda))
}
save(AccSI,file="output/multi_timepoint/accuracies.RData")

```

**7.6.3 Optimal multi vs single time-point selection indices**

As for the single time-point case, one can obtain an optimal regularized selection index such that the accuracy of selection in the testing set is maximum. Code in **Box 10** below can be used to select the optimal PC-SI and L1-PSI, and to compare the accuracy of selection with that of the single time-point CSI, PC-SI and L1-PSI. Results are reported as an average across all partitions along with a 95% confidence interval.

**Box 10.** Optimal selection index

```

tp <- 9      # Time-point
load(paste0("output/timepoint_",tp,"/accuracies.RData"))
load("output/multi_timepoint/accuracies.RData")

AccSI <- rbind(accSI,AccSI)

AccSI$SI <- unlist(lapply(strsplit(as.character(AccSI$SI),"\\"),function(x)x[1]))
AccSI <- split(AccSI,paste(AccSI$rep,"_",AccSI$SI))
AccSI <- do.call(rbind,lapply(AccSI,function(x)x[which.max(x$accuracy),]))

dat = aggregate(accuracy ~ SI, mean,data=AccSI)
dat$sd = aggregate(accuracy ~ SI, sd, data=AccSI)$accuracy
dat$n = aggregate(accuracy ~ SI, length, data=AccSI)$accuracy
dat$se = qnorm(0.975)*dat$sd/sqrt(dat$n)

ggplot(dat,aes(SI,accuracy,ymin=accuracy-se,ymax=accuracy+se)) + theme_bw() +
  geom_bar(stat="identity",width=0.65,fill="orange") + geom_errorbar(width=0.2) +
  geom_text(aes(label=sprintf("%.3f",accuracy)),y=min(dat$accuracy)*0.75+labs(x=""))

```

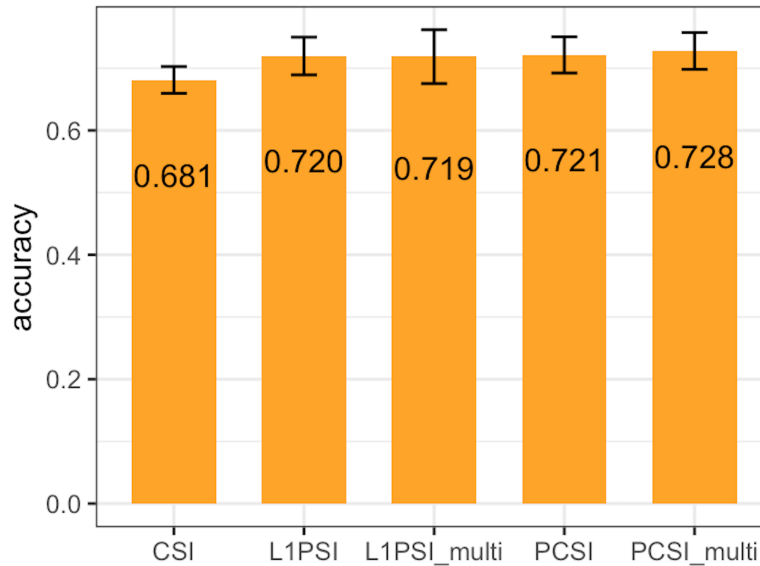


Figure 4: Accuracy of selection (average across 10 TRN-TST partitions) of the CSI, and the optimal PC-SI and L1-PSI for the time-point 9 and multi time-point. Environment E1

## 7.7 Sparsity of the index

The value of  $\lambda$  in (6) for the L1-PSI define the number of predictors with non-zero coefficient in the index. Code in **Box 11** below can be used to estimate the regression coefficients for a multi time-point L1-PSI for the whole dataset. To achieve this, the genetic covariances will be calculated for the 2,250 wavelengths using whole data by means of the `getGenCov` function. The coefficients are then obtained using the `solveEN` with a value of  $\lambda$  obtained from the partitions.

A heatmap of the coefficients will be generated to show the sparsity and regions of the electromagnetic spectrum that were selected by the method.

**Box 11. Sparsity of the optimal index**

```

load("output/varComps.RData")
load("output/multi_timepoint/accuracies.RData")

# Get genetic covariances
x <- scale(do.call(cbind,WL))
fm <- getGenCov(y1=Y$YLD,y2=x,U=evdZZt$vectors,d=evdZZt$values, mc.cores=5)
gencov <- fm$covU

# Get a value of lambda across partitions
AccSI$SI <- unlist(lapply(strsplit(as.character(AccSI$SI),"\\"),function(x)x[1]))
AccSI <- split(AccSI,paste0(AccSI$SI,"_",AccSI$rep))
AccSI <- do.call(rbind,lapply(AccSI,function(x)x[which.max(x$accuracy),]))
lambda <- mean(AccSI$lambda[AccSI$SI=="L1PSI_multi"])

# Get regression coefficients
VARx <- var(x)
beta <- solveEN(VARx,gencov,lambda=lambda)$beta

wl <- factor(rep(gsub("wl","",colnames(WL[[1]])),9))
Timepoint <- factor(rep(1:9,each=ncol(WL[[1]])))
dat <- data.frame(beta=as.vector(beta),Timepoint,wl)
dat$beta[abs(dat$beta) < .Machine$double.eps] <- NA

ggplot(dat, aes(x=wl,y=Timepoint,fill = abs(beta))) + geom_tile(height=0.8) +
  theme_bw() + labs(x="Wavelength (nm)") +
  scale_x_discrete(breaks = levels(wl)[seq(1,nlevels(wl),25)]) +
  scale_fill_gradientn(na.value = 'gray35', colours = c(rev(heat.colors(15))))

```

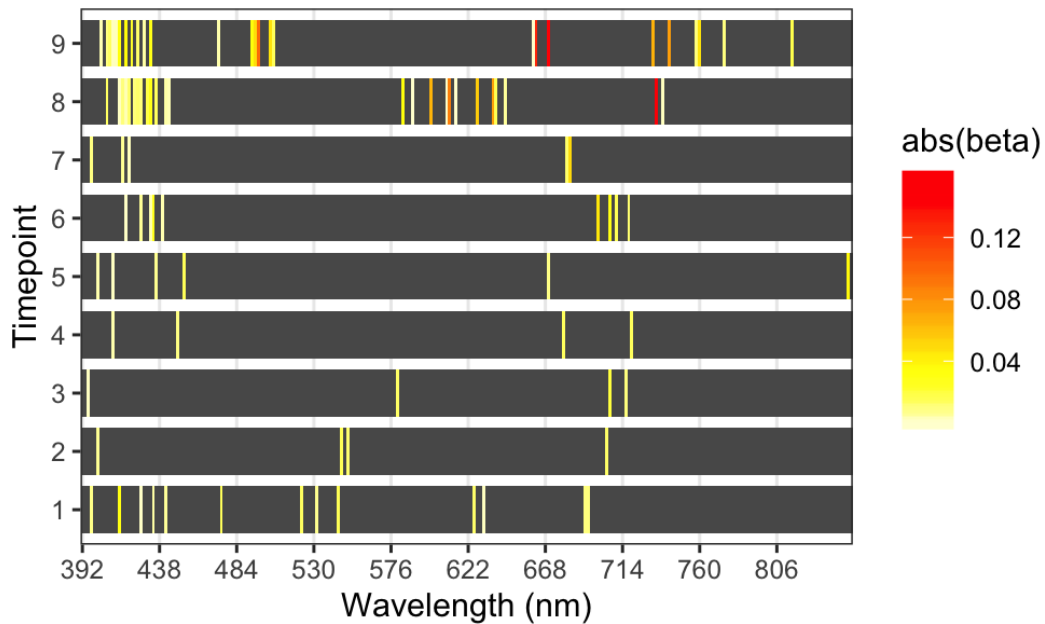


Figure 5: Heatmap of regression coefficients for L1-penalized SI using multi time-point data



## References

- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004). Least angle regression. *The Annals of Statistics* 32(2), 407–499.
- Friedman, J., T. Hastie, H. Höfling, and R. Tibshirani (2007). Pathwise coordinate optimization. *The Annals of Applied Statistics* 1(2), 302–332.
- Fu, W. J. (1998). Penalized regressions: the Bridge versus the LASSO. *Journal of Computational and Graphical Statistics* 7(3), 397–416.
- Hastie, T., R. Tibshirani, and J. H. Friedman (2009). *The elements of statistical learning: data mining, inference, and prediction* (2nd ed.). New York, USA: Springer.
- Hazel, L. N. (1943). The genetic basis for constructing selection indexes. *Genetics* 28(6), 476–490.
- Hoerl, A. E. and R. W. Kennard (1970). Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics* 12(1), 55–67.
- Smith, H. F. (1936). A discriminant function for plant selection. *Annals of Eugenics* 7, 240–250.
- Tibshirani, R. (1996). Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society B* 58(1), 267–288.
- Zou, H. and T. Hastie (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society B* 67(2), 301–320.