

Penalized Functional Binary Regression: An Application to PBC Dataset

Minjie Fan

Instructor: Hans-Georg Müller

March 19, 2013

1 Introduction

Functional binary regression models have proven useful for longitudinal data, in which predictors are usually sparse and irregular. Müller (2005) applied the extension of functional binary regression models to the data called primary biliary cirrhosis(PBC). The data comes from a Mayo Clinic trial conducted from 1974 to 1984. There are various sparsely and irregularly sampled covariates, and survival or censoring time for 312 patients in the data. The aim is to predict long-term survival based on a series of sparse initial measurements. In his paper, Müller merely considered serum bilirubin measurements as predictor, based on which he got the misclassification rate 26.54%. Actually, the other predictors contained in the data are also closely related with the disease, such as albumin and prothrombin time. Therefore, it is necessary to revisit this data and consider all the predictors in the model, for the purpose of decreasing the misclassification error.

Variables in the dataset, see Appendix A.

2 Functional data analysis framework

2.1 Generalized functional linear model

Assume n i.i.d. observations are $(Z_i^1, \dots, Z_i^p, \{X_i^1(t), t \in \mathcal{T}_1\}, \dots, \{X_i^q(t), t \in \mathcal{T}_q\}, Y_i), i = 1, \dots, n$, where Z_i^j 's are random predictor variables, $X_i^j(t)$'s are random predictor curves and Y_i 's are random dependent variables. $E(X_i^j(t)) = \mu^j(t)$. Suppose there is a link function $g(\cdot)$ and a variance function $\sigma^2(\cdot)$. Then a generalized functional linear model or functional quasi-likelihood model is defined as:

$$E(Y_i|Z_i^j, X_i^k(t), j = 1, \dots, p, k = 1, \dots, q) = \mu_i,$$

$$Var(Y_i|Z_i^j, X_i^k(t), j = 1, \dots, p, k = 1, \dots, q) = \sigma^2(\mu_i),$$

$$\begin{aligned}
g(\mu_i) &= \eta_i = \alpha + \beta_1 Z_i^1 + \cdots + \beta_p Z_i^p + \int_{\mathcal{T}_1} (X_i^1(t) - \mu^1(t)) \beta_1(t) dt + \cdots + \int_{\mathcal{T}_q} (X_i^q(t) - \mu^q(t)) \beta_q(t) dt, \\
Y_i &= g^{-1} \left(\alpha + \beta_1 Z_i^1 + \cdots + \beta_p Z_i^p + \int_{\mathcal{T}_1} (X_i^1(t) - \mu^1(t)) \beta_1(t) dt + \cdots + \int_{\mathcal{T}_q} (X_i^q(t) - \mu^q(t)) \beta_q(t) dt \right) + \epsilon_i.
\end{aligned} \tag{1}$$

This model is an extension of the one proposed by Müller and Stadtmüller (2005), which merely contains one random predictor curve $X(t)$.

2.2 Functional principal component analysis

Assume the random curve $X_i^j(t) \in L^2(\mathcal{T}_j)$, abbreviated as $X(t)$ hereinafter. $\mu(t) = E(X(t))$, and $G(s, t) = Cov\{X(s), X(t)\}$. Define the auto-covariance operator

$$(A_G f)(t) = \int f(s) G(s, t) ds,$$

and assume its orthonormal eigenfunctions Φ_k and ordered eigenvalues λ_k exist. Then $G(s, t) = \sum_{k=1}^{\infty} \lambda_k \Phi_k(s) \Phi_k(t)$. By the K-L expansion,

$$X(t) = \mu(t) + \sum_{k=1}^{\infty} \xi_k \Phi_k(t), \tag{2}$$

where ξ_k 's are uncorrelated random variables, called functional principal component scores, and

$$\xi_k = \int (X(t) - \mu(t)) \Phi_k(t) dt.$$

Suppose $\beta_j(t)$ in (1) can be expanded as $\beta_j(t) = \sum_{k=1}^{\infty} \beta_k^j \Phi_k(t)$. Then, the linear predictor in (1) is simplified as

$$\eta_i = \alpha + \beta_1 Z_i^1 + \cdots + \beta_p Z_i^p + \sum_{k=1}^{\infty} \beta_k^1 \xi_k^1 + \cdots + \sum_{k=1}^{\infty} \beta_k^q \xi_k^q. \tag{3}$$

When dealing with sparse and irregular functional data, we use the nonparametric approach called Principal Analysis through Conditional Expectation (PACE), which is under the Gaussian assumptions, see Yao (2005) for details.

3 Fitting method

3.1 Selecting the number of functional principal components

For practical application, we truncate the infinite terms in (3) at $K_j, j = 1, \dots, q$,

$$\eta_i = \alpha + \beta_1 Z_i^1 + \cdots + \beta_p Z_i^p + \sum_{k=1}^{K_1} \beta_k^1 \xi_k^1 + \cdots + \sum_{k=1}^{K_q} \beta_k^q \xi_k^q. \tag{4}$$

Define the fraction of variance explained (FVE) as $F(K) = \frac{\sum_{k=1}^K \lambda_k}{\sum_{k=1}^{\infty} \lambda_k}$. We select the smallest K such that $F(K)$ is larger than a pre-specified threshold, which is taken as 0.99.

3.2 Predictor selection by group lasso

The number of predictors, except the intercept, is $n = p + \sum_{i=1}^q K_i$. when n is large, these predictors compose a high-dimensional dataset. Therefore, it is necessary to do the predictor selection. Besides, the predictors from the same random predictor curve should be selected simultaneously.

Meier (2008) proposed the group lasso for logistic regression. Suppose the vector of parameters is $(\beta_0, \boldsymbol{\beta}^T)^T$, with dimension $p + 1$, where β_0 is the intercept. $\boldsymbol{\beta}$ is grouped in G groups $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_G$, with group sizes $\{\text{df}_1, \dots, \text{df}_G\}$. The logistic group lasso estimator $\hat{\boldsymbol{\beta}}_\lambda$ is given by the minimizer of the convex function

$$S_\lambda(\boldsymbol{\beta}) = -l(\boldsymbol{\beta}) + \lambda \sum_{g=1}^G s(\text{df}_g) \|\boldsymbol{\beta}_g\|_2, \quad (5)$$

where $l(\cdot)$ is the log-likelihood function, df_g is the number of parameters in the g^{th} group, and the function $s(\text{df})$ is used to rescale the penalty, which is taken as $\text{df}_g^{1/2}$.

The value of λ can be chosen by k -fold Cross-validation, where k is taken as 10.

The design matrix $X = [X_1 | \dots | X_G]$ should be centralized and block-wise standardized before entering the model, where we assume X_g 's are of full column rank.

4 Application to PBC dataset

4.1 Preprocessing

We include the patients who satisfy the following two conditions:

- Survived the first 910 days of the study;
- Survival status beyond 10 years was known.

Altogether, there are 260 patients satisfied these two conditions, of which 84 died between 910 and 3650 days, and 176 lived beyond 3650 days. Based on the several initial measurements during the first 910 days, we would like to predicting survival beyond 10 years after entering the study. Omitting the predictors who have missing values or are categorical, we select 7 predictors, which are drug, age, sex, serum bilirubin, albumin, prothrombin time (PT) and SGOT. Among them, drug, age and sex are random variables, while the remaining are random curves. The bilirubin and SGOT measurements are log-transformed. For the responses Y_i ,

we use 0 to represent short-lived, and 1 to represent long-lived.

The outliers are removed w.r.t. the random curves one by one. We find that there are outliers in albumin and PT, see Appendix B, Figure 1 and 2. 4 short-lived and 3 long-lived patients are removed in total.

4.2 Model fitting

By PACE, we get the estimates of mean functions, covariance surface, and eigenfunctions for each of the four random predictor curves, see Appendix B, Figure 3-7. We find that the $\log(\text{bilirubin})$, PT and $\log(\text{SGOT})$ are much higher for the short-lived patients than for the long-lived ones; the albumin is much lower for the short-lived patients. Besides, the $\log(\text{bilirubin})$, PT and $\log(\text{SGOT})$ are increasing over time, and the albumin is decreasing over time, for short-lived patients. These results are consistent with the ones of clinic study. Therefore, the four random predictor curves are good bio-markers for PBC patients, which distinguish short-lived patients from long-lived ones.

By FVE criterion, serum bilirubin, albumin, PT and SGOT have 3,4,4 and 1 eigenfunctions after truncation, respectively. Combined with the 3 random predictor variables drug, age and sex, we get the design matrix X , which is a 253×15 matrix (doesn't include the intercept term).

We use the package `grplasso` in R software to realize the method of group lasso. The multiplicative grid search set of the penalty parameter λ is $\{\lambda_{\max}, 0.96\lambda_{\max}, \dots, 0.96^{148}\lambda_{\max}, 0\}$. λ is selected by 10-fold Cross-validation.

4.3 Results

Table 1: Leave-one-out classification detailed results

	Long-lived (True)	Short-lived (True)
Long-lived (Classified)	156	40
Short-lived (Classified)	17	40

Using the leave-one-out prediction error criterion, the overall misclassification rate is 22.53%, which is smaller than the one in Müller (2005), 26.54%. By including more predictors and using the shrinkage method, we decrease the misclassification error efficiently. Table 1 shows the detailed results. The misclassification rate for the short-lived patients are 50%, while the one for the long-lived patients are 9.83%. This can be explained by the unbalanced sample sizes of the short-lived and long-lived categories.

The histogram of the optimal penalty parameter λ selected by Cross-validation in the 253-times leave-one-out fittings is shown in Appendix B, Figure 8. The plot is right-skewed

Table 2: Parameter estimates by group lasso

Predictors	Coefficients
Intercept	4.5510001249
1 st score of Bili	-0.0207121621
2 nd score of Bili	-0.0125968522
3 rd score of Bili	0.0084731104
1 st score of Alb	0.0674235645
2 nd score of Alb	-0.1327425617
3 rd score of Alb	0.1218763442
4 th score of Alb	-1.7952969754
1 st score of PT	0.0004362802
2 nd score of PT	0.0550491950
3 rd score of PT	0.2886511832
4 th score of PT	-0.0769774141
1 st score of SGOT	-0.0325174314
drug	-0.1434167985
age	-0.0708683351
sex	0.0742933828

Table 3: The estimated importance degrees of all the predictors in decreasing order

serum bilirubin
age
prothrombin time
SGOT
albumin
sex
drug

and most λ 's are between 0 and 0.25. Fitting the model by all the observations, we get the optimal λ , which is 0.1681001, by Cross-validation. The fitted coefficients are listed in table 2. The estimated standard errors of these fitted coefficients can be obtained by bootstrapping. The interpretation of the signs of these fitted coefficients is not reliable because the primary aim of the group lasso is to select the optimal set of predictors, not to find the predictors that give best explanation. However, taking a look at the access order of all the predictors into the model with λ varying from λ_{\max} to 0, we can rate their importance degrees according to the order. Table 3 shows the estimated importance degrees of all the predictors in decreasing order. This makes sense because serum bilirubin concentration is one of the most important indicators of chronic liver cirrhosis, such as PBC, and the survival time is closely related with the age of the patient.

4.4 Discussion

Functional data analysis provides an innovative approach to handle random trajectories and infinite-dimensional data. Like, the classification of longitudinal data can be realized by functional binary regression, even in the case of sparse and irregular predictors. However, the problem of predictor selection still exists, especially for the data that have thousands or millions of random predictor variables and random predictor curves, such as fMRI images and gene sequences. Combined shrinkage methods with functional data analysis is a possible solution to this problem.

References

- [1] Meier, L., Van De Geer, S. and Bühlmann, P. 2008. The group lasso for logistic regression. *Journal of the Royal Statistical Society: Series B*, **70**, 53–71.
- [2] Müller, H.-G. and Stadtmüller, U. 2005. Generalized functional linear models. *Annals of statistics*, **33**, 774–805.
- [3] Müller, H.-G. 2005. Functional modelling and classification of longitudinal data. *Scandinavian Journal of Statistics*, **32**, 223–240.
- [4] Yao, F., Müller, H.-G. and Wang, J. 2005. Functional data analysis for sparse longitudinal data. *Journal of the American Statistical Association*, **100**, 577–590.

5 Appendix

5.1 Appendix A

Variables:

case number

number of days between registration and the earlier of death, transplantation, or study analysis time

status: 0=alive, 1=transplanted, 2=dead

drug: 1=D-penicillamine, 0=placebo

age in days, at registration

sex: 0=male, 1=female

day: number of days between enrollment and this visit date, remaining values on the line of data refer to this visit.

presence of ascites: 0=no 1=yes

presence of hepatomegaly: 0=no 1=yes

presence of spiders: 0=no 1=yes

presence of edema: 0=no edema and no diuretic therapy for edema; .5=edema present without diuretics, or edema resolved by diuretics; 1=edema despite diuretic therapy

serum bilirubin in mg/dl

serum cholesterol in mg/dl

albumin in gm/dl

alkaline phosphatase in U/liter

SGOT in U/ml (serum glutamic-oxaloacetic transaminase, the enzyme name has subsequently changed to ALT in the medical literature)

platelets per cubic in ml/1000

prothrombin time in seconds

histologic stage of disease

5.2 Appendix B

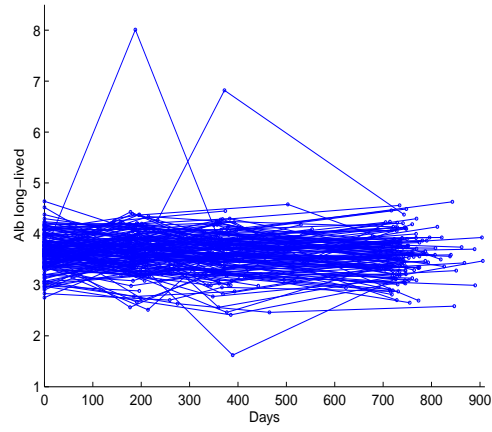


Figure 1: The observed initial albumin measurements, for 176 long-lived patients.

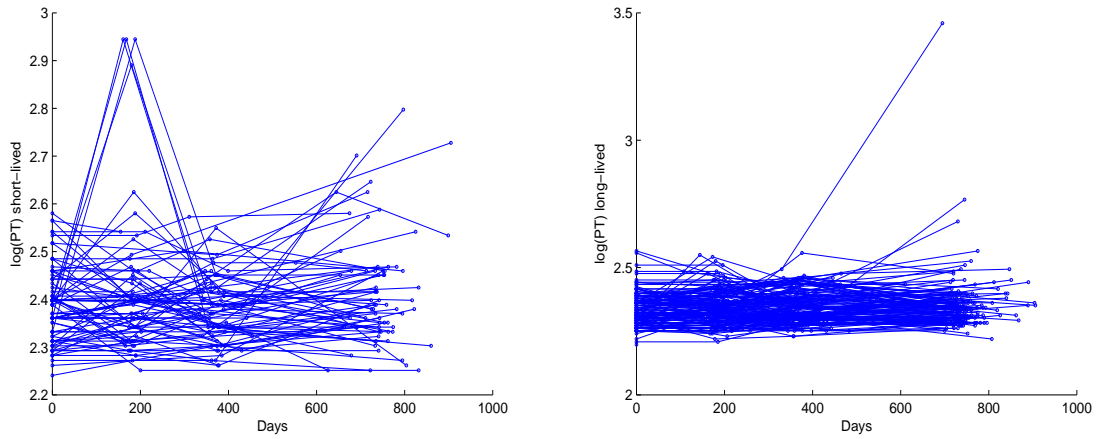


Figure 2: The observed initial log(PT) measurements (Left: for 84 short-lived patients, Right: for 174 long-lived patients).

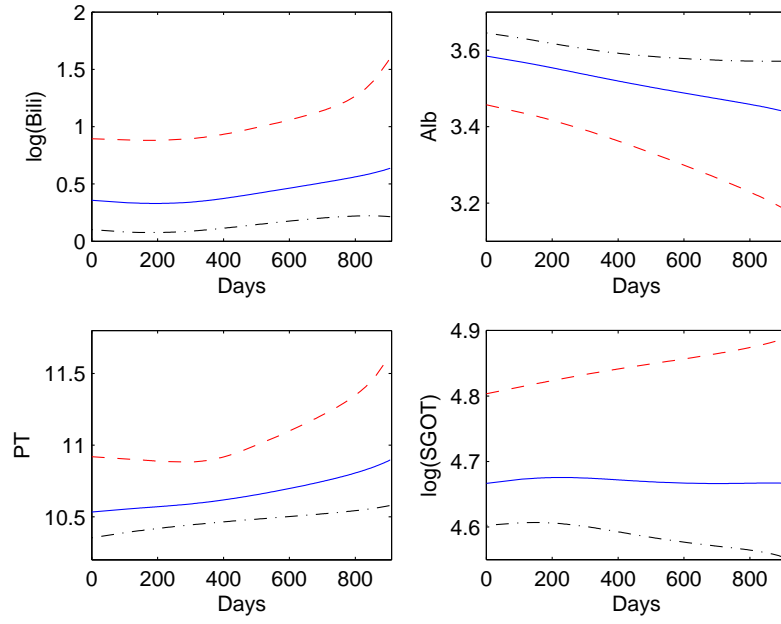


Figure 3: Smoothed mean functions for four random predictor curves (solid: all the 253 patients, dashed: the short-lived patients, dotdashed: the long-lived patients).

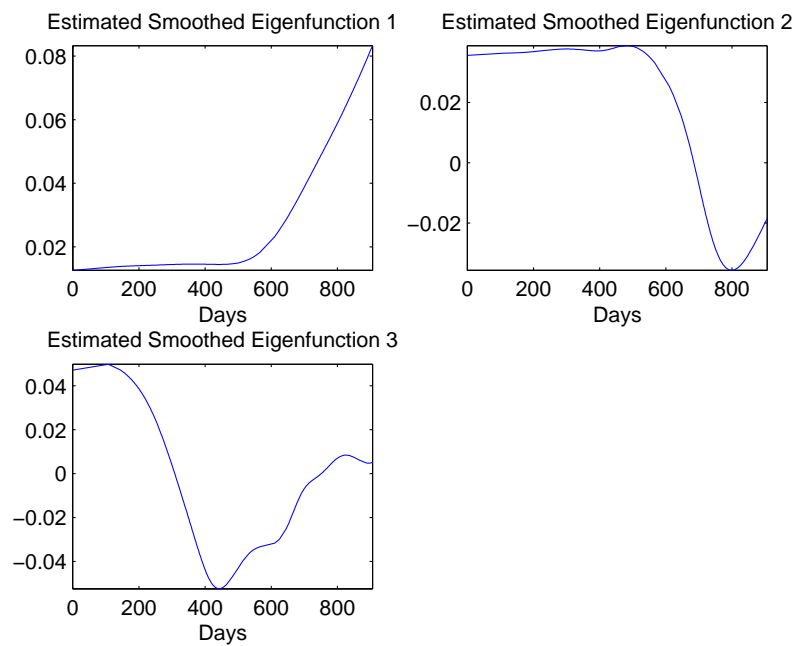


Figure 4: Smoothed eigenfunctions for serum bilirubin.

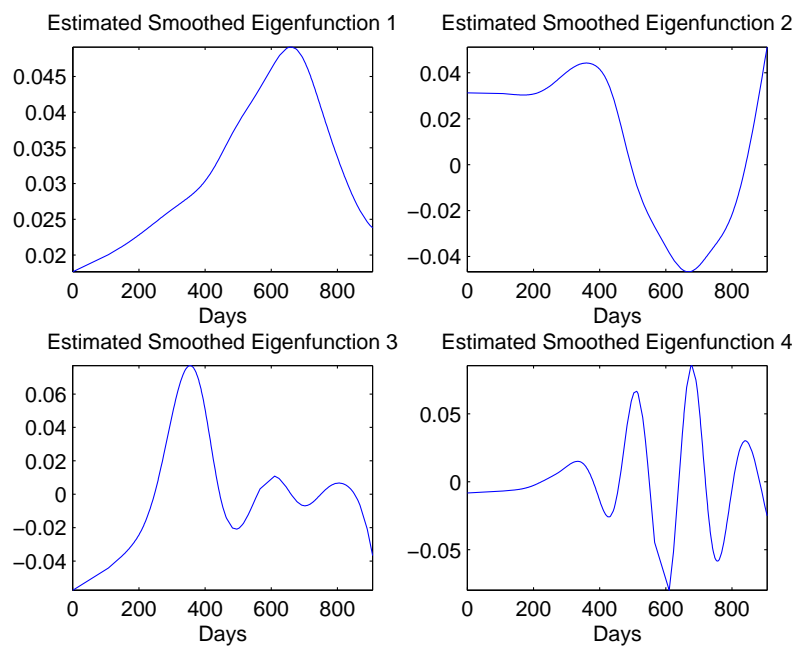


Figure 5: Smoothed eigenfunctions for albumin.

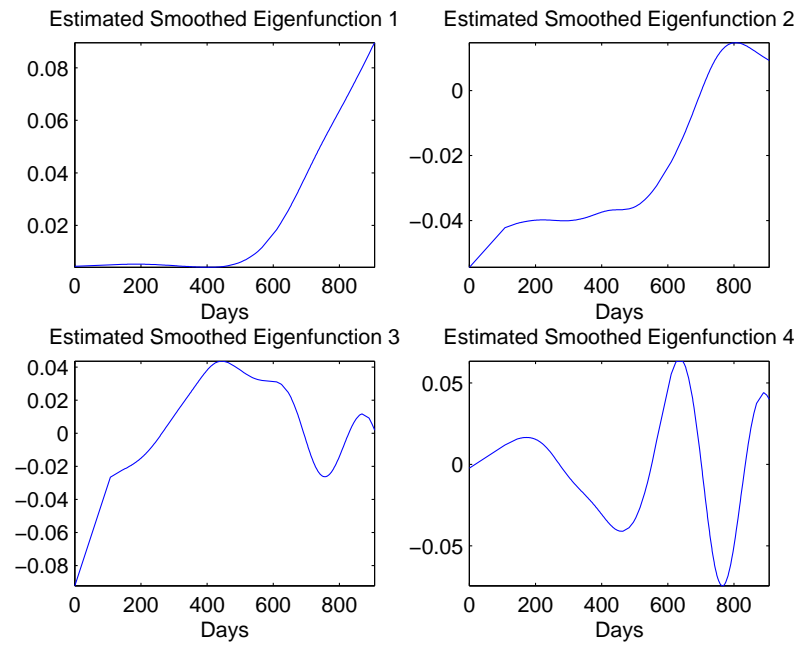


Figure 6: Smoothed eigenfunctions for PT.

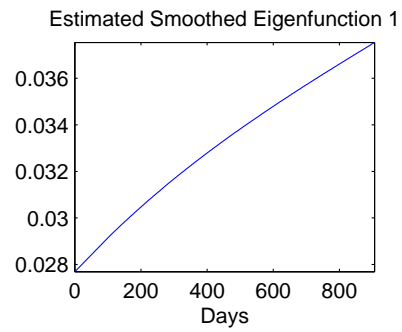


Figure 7: Smoothed eigenfunctions for SGOT.

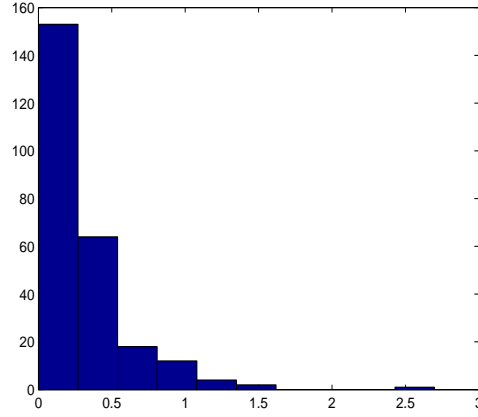


Figure 8: The histogram of the optimal penalty parameter λ selected by Cross-validation in the 253-times leave-one-out fittings.

5.3 Appendix C

Main.m (Matlab)

```

1 clear
2 clc
3 load('PBC.mat')
4
5 %pre-processing
6 PBC_910=PBC((PBC(:,2) >= 910),:);
7 PBC_died=PBC_910((PBC_910(:,3)==2).*(PBC_910(:,2) <= 3650)==1,:);
8 index_died=unique(PBC_died(:,1));
9 PBC_live=PBC_910((PBC_910(:,3)~=2)+(PBC_910(:,2) > 3650) >= 1,:);
10 index_live=unique(PBC_live(:,1));
11
12 PBC_died_910=PBC_died(PBC_died(:,7) <= 910,:);
13 PBC_live_910=PBC_live(PBC_live(:,7) <= 910,:);
14
15 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
16
17 %plot log(bili) w.r.t short-lived
18 subplot(2,1,1)
19 hold on
20 for i=1:length(index_died)
21     index_sub=index_died(i);
22     PBC_sub=PBC_died_910(PBC_died_910(:,1)==index_sub,:);
23     plot(PBC_sub(:,7), log(PBC_sub(:,12)), '-o', 'MarkerSize', 2)
24 end
25 axis([0 910 -3 4])
26 xlabel('Days')
27 ylabel('log(Bili) short-lived')
28
29 %plot log(bili) w.r.t long-lived
30 subplot(2,1,2)
31 hold on
32 for i=1:length(index_live)
33     index_sub=index_live(i);
34     PBC_sub=PBC_live_910(PBC_live_910(:,1)==index_sub,:);
35     plot(PBC_sub(:,7), log(PBC_sub(:,12)), '-o', 'MarkerSize', 2)
36 end
37 axis([0 910 -3 4])
38 xlabel('Days')
39 ylabel('log(Bili) long-lived')
40
41 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
42

```

```

43 %plot alb w.r.t short-lived
44 subplot(2,1,1)
45 hold on
46 for i=1:length(index_died)
47     index_sub=index_died(i);
48     PBC_sub=PBC_died_910(PBC_died_910(:,1)==index_sub,:);
49     plot(PBC_sub(:,7),PBC_sub(:,14),'-o','MarkerSize',2)
50 end
51 axis([0 910 1.5 5])
52 xlabel('Days')
53 ylabel('Alb short-lived')
54
55 %remove the outlier
56 index_live=remove_outlier(PBC_live_910,6.5,index_live,14,false);
57
58 %plot alb w.r.t long-lived
59 subplot(2,1,2)
60 hold on
61 for i=1:length(index_live)
62     index_sub=index_live(i);
63     PBC_sub=PBC_live_910(PBC_live_910(:,1)==index_sub,:);
64     plot(PBC_sub(:,7),PBC_sub(:,14),'-o','MarkerSize',2)
65 end
66 axis([0 910 1.5 5])
67 xlabel('Days')
68 ylabel('Alb long-lived')
69
70 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
71
72 %remove the outlier
73 index_died=remove_outlier(PBC_died_910,2.85,index_died,18,true);
74
75 %plot PT w.r.t short-lived
76 subplot(2,1,1)
77 hold on
78 for i=1:length(index_died)
79     index_sub=index_died(i);
80     PBC_sub=PBC_died_910(PBC_died_910(:,1)==index_sub,:);
81     plot(PBC_sub(:,7),PBC_sub(:,18),'-o','MarkerSize',2)
82 end
83 axis([0 910 8 18])
84 xlabel('Days')
85 ylabel('PT short-lived')
86
87 %remove the outlier
88 index_live=remove_outlier(PBC_live_910,3,index_live,18,true);
89
90 %plot PT w.r.t long-lived
91 subplot(2,1,2)
92 hold on
93 for i=1:length(index_live)
94     index_sub=index_live(i);
95     PBC_sub=PBC_live_910(PBC_live_910(:,1)==index_sub,:);
96     plot(PBC_sub(:,7),PBC_sub(:,18),'-o','MarkerSize',2)
97 end
98 axis([0 910 8 18])
99 xlabel('Days')
100 ylabel('PT long-lived')
101
102 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
103
104 %plot log(SGOT) w.r.t short-lived
105 subplot(2,1,1)
106 hold on
107 for i=1:length(index_died)
108     index_sub=index_died(i);
109     PBC_sub=PBC_died_910(PBC_died_910(:,1)==index_sub,:);
110     plot(PBC_sub(:,7),log(PBC_sub(:,16)),'-o','MarkerSize',2)
111 end
112 axis([0 910 3 7])
113 xlabel('Days')
114 ylabel('log(SGOT) short-lived')
115
116 %plot log(SGOT) w.r.t long-lived
117 subplot(2,1,2)
118 hold on
119 for i=1:length(index_live)

```

```

120     index_sub=index_live(i);
121     PBC_sub=PBC_live_910(PBC_live_910(:,1)==index_sub,:);
122     plot(PBC_sub(:,7),log(PBC_sub(:,16)),'-o','MarkerSize',2)
123 end
124 axis([0 910 3 7])
125 xlabel('Days')
126 ylabel('log(SGOT) long-lived')
127
128 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
129
130 %preliminary computation
131 n_died=length(index_died);
132 n_live=length(index_live);
133 n_total=n_died+n_live;
134
135 %set
136 p=setOptions('selection_k','FVE','FVE_threshold',0.99);
137
138 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
139
140 %FPCA(Bili short-lived)
141 [X_Bili_died,y_died,t_died]=my_FPCA(PBC_died_910,n_died,index_died,12,true);
142 out1_Bili_died=getVal(X_Bili_died,'out1');
143 mu_Bili_died=getVal(X_Bili_died,'mu');
144
145 %FPCA(Bili long-lived)
146 [X_Bili_live,y_live,t_live]=my_FPCA(PBC_live_910,n_live,index_live,12,true);
147 out1_Bili_live=getVal(X_Bili_live,'out1');
148 mu_Bili_live=getVal(X_Bili_live,'mu');
149
150 %FPCA(Bili overall)
151 y_tot=[y_died y_live];
152 t_tot=[t_died t_live];
153 X_Bili_tot=FPCA(y_tot,t_tot,p); %output
154 out1_Bili_tot=getVal(X_Bili_tot,'out1');
155 mu_Bili_tot=getVal(X_Bili_tot,'mu');
156
157 %plot
158 subplot(2,2,1)
159 plot(out1_Bili_died,mu_Bili_died,'-r')
160 hold on
161 plot(out1_Bili_live,mu_Bili_live,'-k')
162 plot(out1_Bili_tot,mu_Bili_tot,'b')
163 axis([0 910 0 2])
164 xlabel('Days')
165 ylabel('log(Bili)')
166
167 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
168
169 %FPCA(Alb short-lived)
170 [X_Alb_died,y_died,t_died]=my_FPCA(PBC_died_910,n_died,index_died,14,false);
171 out1_Alb_died=getVal(X_Alb_died,'out1');
172 mu_Alb_died=getVal(X_Alb_died,'mu');
173
174 %FPCA(Alb long-lived)
175 [X_Alb_live,y_live,t_live]=my_FPCA(PBC_live_910,n_live,index_live,14,false);
176 out1_Alb_live=getVal(X_Alb_live,'out1');
177 mu_Alb_live=getVal(X_Alb_live,'mu');
178
179 %FPCA(Alb overall)
180 y_tot=[y_died y_live];
181 t_tot=[t_died t_live];
182 X_Alb_tot=FPCA(y_tot,t_tot,p); %output
183 out1_Alb_tot=getVal(X_Alb_tot,'out1');
184 mu_Alb_tot=getVal(X_Alb_tot,'mu');
185
186 %plot
187 subplot(2,2,2)
188 plot(out1_Alb_died,mu_Alb_died,'-r')
189 hold on
190 plot(out1_Alb_live,mu_Alb_live,'-k')
191 plot(out1_Alb_tot,mu_Alb_tot,'b')
192 axis([0 910 3.1 3.7])
193 xlabel('Days')
194 ylabel('Alb')
195
196 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

```

```

197
198 %FPCA(PT short-lived)
199 [X_PT_died, y_died, t_died]=my_FPCA(PBC_died_910, n_died, index_died, 18, false);
200 out1_PT_died=getVal(X_PT_died, 'out1');
201 mu_PT_died=getVal(X_PT_died, 'mu');
202
203 %FPCA(PT long-lived)
204 [X_PT_live, y_live, t_live]=my_FPCA(PBC_live_910, n_live, index_live, 18, false);
205 out1_PT_live=getVal(X_PT_live, 'out1');
206 mu_PT_live=getVal(X_PT_live, 'mu');
207
208 %FPCA(PT overall)
209 y_tot=[y_died y_live];
210 t_tot=[t_died t_live];
211 X_PT_tot=FPCA(y_tot, t_tot, p); %output
212 out1_PT_tot=getVal(X_PT_tot, 'out1');
213 mu_PT_tot=getVal(X_PT_tot, 'mu');
214 getVal(X_PT_tot, 'no_opt')
215
216 %plot
217 subplot(2,2,3)
218 plot(out1_PT_died, mu_PT_died, '--r')
219 hold on
220 plot(out1_PT_live, mu_PT_live, '--k')
221 plot(out1_PT_tot, mu_PT_tot, 'b')
222 axis([0 910 10.2 11.8])
223 xlabel('Days')
224 ylabel('PT')
225
226 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
227
228 %FPCA(SGOT short-lived)
229 [X_SGOT_died, y_died, t_died]=my_FPCA(PBC_died_910, n_died, index_died, 16, true);
230 out1_SGOT_died=getVal(X_SGOT_died, 'out1');
231 mu_SGOT_died=getVal(X_SGOT_died, 'mu');
232
233 %FPCA(SGOT long-lived)
234 [X_SGOT_live, y_live, t_live]=my_FPCA(PBC_live_910, n_live, index_live, 16, true);
235 out1_SGOT_live=getVal(X_SGOT_live, 'out1');
236 mu_SGOT_live=getVal(X_SGOT_live, 'mu');
237
238 %FPCA(SGOT overall)
239 y_tot=[y_died y_live];
240 t_tot=[t_died t_live];
241 X_SGOT_tot=FPCA(y_tot, t_tot, p); %output
242 out1_SGOT_tot=getVal(X_SGOT_tot, 'out1');
243 mu_SGOT_tot=getVal(X_SGOT_tot, 'mu');
244 getVal(X_SGOT_tot, 'no_opt')
245
246 %plot
247 subplot(2,2,4)
248 plot(out1_SGOT_died, mu_SGOT_died, '--r')
249 hold on
250 plot(out1_SGOT_live, mu_SGOT_live, '--k')
251 plot(out1_SGOT_tot, mu_SGOT_tot, 'b')
252 axis([0 910 4.55 4.9])
253 xlabel('Days')
254 ylabel('log (SGOT)')
255
256 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
257
258 %plot smoothed eigenfunctions
259 plot_eigenfun(X_Bili_tot)
260 plot_eigenfun(X_Alb_tot)
261 plot_eigenfun(X_PT_tot)
262 plot_eigenfun(X_SGOT_tot)
263
264 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
265
266 %pre-processing
267 PC_scores_Bili=getVal(X_Bili_tot, 'xi_est');
268 PC_scores_Alb=getVal(X_Alb_tot, 'xi_est');
269 PC_scores_PT=getVal(X_PT_tot, 'xi_est');
270 PC_scores_SGOT=getVal(X_SGOT_tot, 'xi_est');
271
272 drug=zeros(n_total,1);
273 for i=1:n_died

```

```

274     drug(i,1)=unique(PBC_died_910(find(PBC_died_910(:,1)==index_died(i)),4));
275 end
276 for i=1:n_live
277     drug(n_died+i,1)=unique(PBC_live_910(find(PBC_live_910(:,1)==index_live(i)),4));
278 end
279
280 age=zeros(n_total,1);
281 for i=1:n_died
282     age(i,1)=unique(PBC_died_910(find(PBC_died_910(:,1)==index_died(i)),5))/365;
283 end
284 for i=1:n_live
285     age(n_died+i,1)=unique(PBC_live_910(find(PBC_live_910(:,1)==index_live(i)),5))/365;
286 end
287
288 sex=zeros(n_total,1);
289 for i=1:n_died
290     sex(i,1)=unique(PBC_died_910(find(PBC_died_910(:,1)==index_died(i)),6));
291 end
292 for i=1:n_live
293     sex(n_died+i,1)=unique(PBC_live_910(find(PBC_live_910(:,1)==index_live(i)),6));
294 end
295
296 label=[zeros(n_died,1);ones(n_live,1)];
297 predictor=[PC_scores_Bili PC_scores_Alb PC_scores_PT PC_scores_SGOT drug age sex];
298 save('glm.mat','label','predictor','-v6')

```

my_FPCA.m (Matlab)

```

1 function [X,y,t]=my_FPCA(PBC_910,n,index,which,logtf)
2 y=cell(1,n);
3 if logtf==true
4     for i=1:n
5         y{i}=log(PBC_910(PBC_910(:,1)==index(i),which));
6     end
7 else
8     for i=1:n
9         y{i}=PBC_910(PBC_910(:,1)==index(i),which);
10    end
11 end
12 t=cell(1,n);
13 for i=1:n
14     t{i}=PBC_910(PBC_910(:,1)==index(i),7);
15 end
16 X=FPCA(y,t);
17 end

```

plot_eigenfun.m (Matlab)

```

1 function plot_eigenfun(X)
2 n=getVal(X,'no_opt');
3 phi=getVal(X,'phi');
4 phi_time=getVal(X,'out1');
5 for i=1:n
6     subplot(2,2,i)
7     plot(phi_time,phi(:,i))
8     xlabel('Days')
9     title(['Estimated Smoothed Eigenfunction ' num2str(i)])
10    axis tight
11 end
12 end

```

remove_outlier.m (Matlab)

```

1 function [index,outlier]=remove_outlier(PBC_910,threshold,index,which,logtf)
2 if logtf==false
3     outlier=unique(PBC_910(PBC_910(:,which)>threshold,1));
4 else
5     outlier=unique(PBC_910(log(PBC_910(:,which))>threshold,1));
6 for i=1:length(outlier)
7     index(find(index==outlier(i)))=[];
8 end
9 end

```


group_lasso.r (R)

```

1  rm(list=ls())
2  library(R.matlab)
3  library(grplasso)
4  library(bootstrap)
5  dt<-readMat("D://glm.mat")
6  label<-dt$label
7  predictor<-dt$predictor
8  n<-nrow(predictor)
9  predictor<-cbind(matrix(rep(1,n),n,1),predictor)
10 index<-c(NA,1,1,1,2,2,2,2,3,3,3,3,4,5,6,7)
11
12 theta.predict<-function(fit.obj,x){
13   predict(fit.obj,x,type="response")
14 }
15
16 cv.error<-function(tra.pre,tra.lab,lambda,k=10){
17   index<-c(NA,1,1,1,2,2,2,2,3,3,3,3,4,5,6,7)
18   list.val<-crossval(tra.pre,tra.lab,grplasso,theta.predict,index=index,
19 model=LogReg(),lambda=lambda,center=TRUE,standardize=TRUE,ngroup=k)
20   mean(as.numeric(list.val$cv.fit>0.5)!=tra.lab)
21 }
22
23 lambda.max<-30
24 m<-148
25 search.set<-rep(0,m+2)
26 search.set[1]<-lambda.max
27 for (i in 2:(m+1)){
28   search.set[i]<-search.set[i-1]*0.96
29 }
30 R<-10
31 pre.lab<-matrix(0,n,1)
32 record.min.lambda<-matrix(0,n,1)
33
34 for (i in 1:n){
35   tra.pre<-predictor[-i,]
36   tra.lab<-label[-i]
37   min.error<-1
38   for (j in 1:(m+2)){
39     lambda<-search.set[j]
40     s<-0
41     for (r in 1:R){
42       s<-s+cv.error(tra.pre,tra.lab,lambda)
43     }
44     s<-s/R
45     if (s<min.error){
46       min.error<-s
47       min.lambda<-lambda
48     }
49   }
50   record.min.lambda[i]<-min.lambda
51   fit.obj<-grplasso(x=tra.pre,y=tra.lab,index=index,model=LogReg(),
52 lambda=min.lambda,center=TRUE,standardize=TRUE)
53   pre.lab[i]<-as.numeric(predict(fit.obj,t(as.matrix(predictor[i,])),type="response")>0.5)
54 }
55 write.csv(pre.lab,"result1.csv")
56 write.csv(record.min.lambda,"result2.csv")
57
58 min.error<-1
59 for (j in 1:(m+2)){
60   lambda<-search.set[j]
61   s<-0
62   for (r in 1:R){
63     s<-s+cv.error(predictor,label,lambda)
64   }
65   s<-s/R
66   if (s<min.error){
67     min.error<-s
68     min.lambda<-lambda
69   }
70 }
71
72 fit.obj<-grplasso(x=predictor,y=label,index=index,model=LogReg(),
73 lambda=min.lambda,center=TRUE,standardize=TRUE)
74 fit.obj2<-glm.fit(x=predictor,y=label)

```