



Feature Selector App

by Vladislav Antipin and Martin Larsen

<https://www.immulab.fr/cms/index.php/team/tools/lab-tools/feature-selector>

Team: Immunity & microbiota ecology

www.immulab.fr

INSERM UMR-S1135, Sorbonne University, Paris, France

1

1

2

Glossary

- Variable = Parameter = Dimension (= Column)
- Feature = Predictor variable
- Response variable = Outcome Variable = variable to be predicted
- Observation = Sample = Point (= Row)

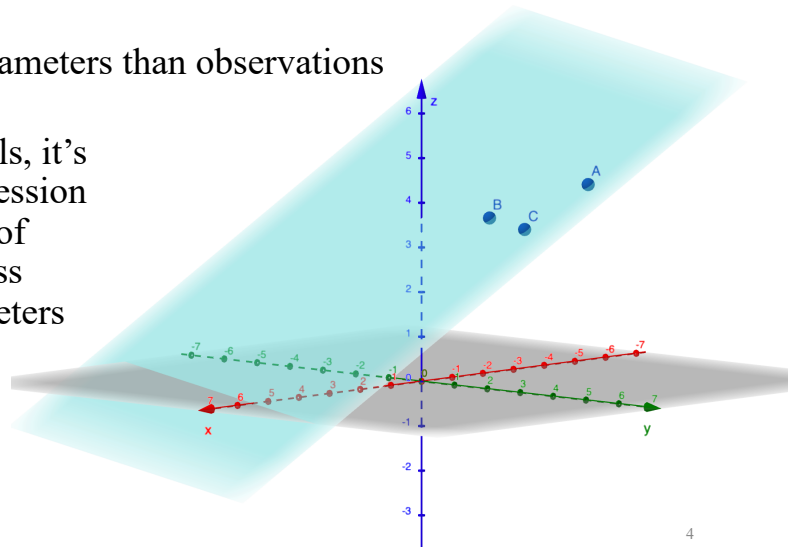
3

3

Aim

$p > n$ problem : more parameters than observations

For most regression models, it's not possible to fit the regression hyperplane if the number of observations (points) is less than the number of parameters (dimensions)



4

4

Aim

Redundancy problem:

we're interested to identify the most important features

IO	IP	IQ	IR	IS	IT	IU	IV	IW	IX	IY	IZ	JA	JB	JC		
nCD4TregT	nCD4Th1	If nCD4	TH17	nCD4	TH22	nCD4N	nCD4cm	nCD4em	nCD4tm	nCD4M	nCD8N	nCD8 Eff	nCD8cm	nCD8em	nCD8tm	nCD8M

5

5

Solution: Regularized models

6

6

Regularized models

Cost Function = Loss Function + Penalty

Lasso (L¹ Norm)

$$\lambda \sum_{j=1}^p |\beta_j|$$

Ridge (L² Norm)

$$\frac{\lambda}{2} \sum_{j=1}^p \beta_j^2$$

7

7

Regularized models

Cost Function = Loss Function + Penalty

Lasso (L¹ Norm)

$$\lambda \sum_{j=1}^p |\beta_j|$$

ElasticNet

$$\lambda \left[\alpha \sum_{j=1}^p |\beta_j| + \frac{(1-\alpha)}{2} \sum_{j=1}^p \beta_j^2 \right]$$

Ridge (L² Norm)

$$\frac{\lambda}{2} \sum_{j=1}^p \beta_j^2$$

8

8

Regularized models

Cost Function = Loss Function + Penalty

Lasso (L¹ Norm)

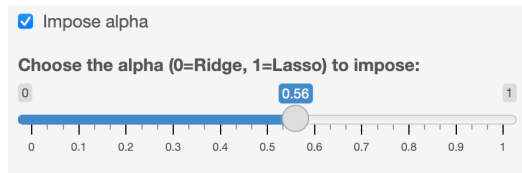
ElasticNet

Ridge (L² Norm)

$$\lambda \sum_{j=1}^p |\beta_j|$$

$$\lambda \left[\alpha \sum_{j=1}^p |\beta_j| + \frac{(1-\alpha)}{2} \sum_{j=1}^p \beta_j^2 \right]$$

$$\frac{\lambda}{2} \sum_{j=1}^p \beta_j^2$$



9

9

Regularized models

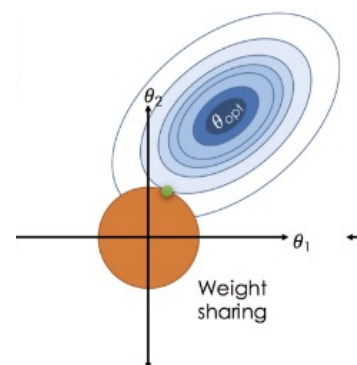
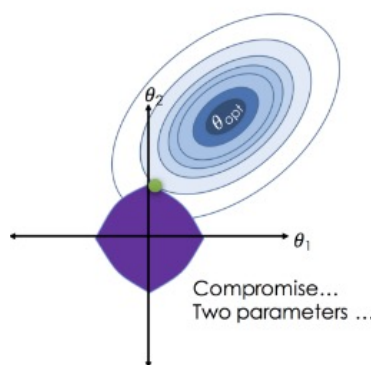
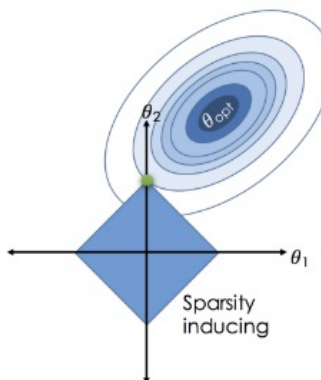
“strict”
“sparse”

Lasso (L¹ Norm)

ElasticNet

Ridge (L² Norm)

“soft”
“grouping”



Zou, H., & Hastie, T. (2005). Regularization and variable selection via the elastic net.

10

Regularized models

Regularized = Penalized (any penalty)
 \approx Sparse (only Lasso)

- regularized GLM (Lasso, Elastic-Net, Ridge) you choose alpha
- sparse PLS-DA (only Lasso)
- Multi-Block sparse PLS-DA = DIABLO (only Lasso)

11

11

Suitable methods

Based on response variable:

Categorical

- binary (e.g. dead-alive) sPLS-DA and GLM
- multinomial (e.g. response - partial response - no response) only sPLS-DA

Numeric (e.g. CRP level in blood)

only GLM

Survival time

- event + time until event only GLM

Based on predictors:

• Numeric + Categorical

GLM

• only Numeric

sPLS-DA

(categorical excluded automatically)

12

Data structure

Data - Monoblock

Features



13

13

Data structure

Data - Monoblock

Features



Feature table

Description



ID variable is obligatory and has to UNIQUELY identify each observation after filters are applied!

check that there are no duplicated rows or columns!

14

14

Data structure

Data - Monoblock

Feature table

Features

Description



15

15

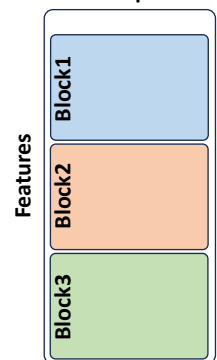
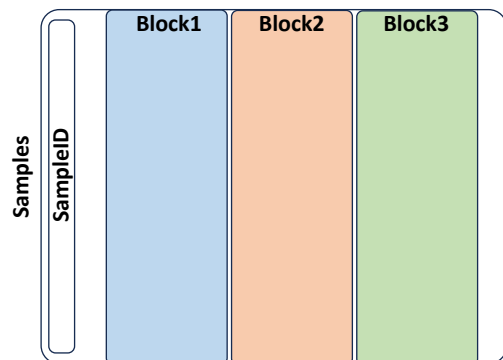
Data structure

Data - Multiblock

Feature table

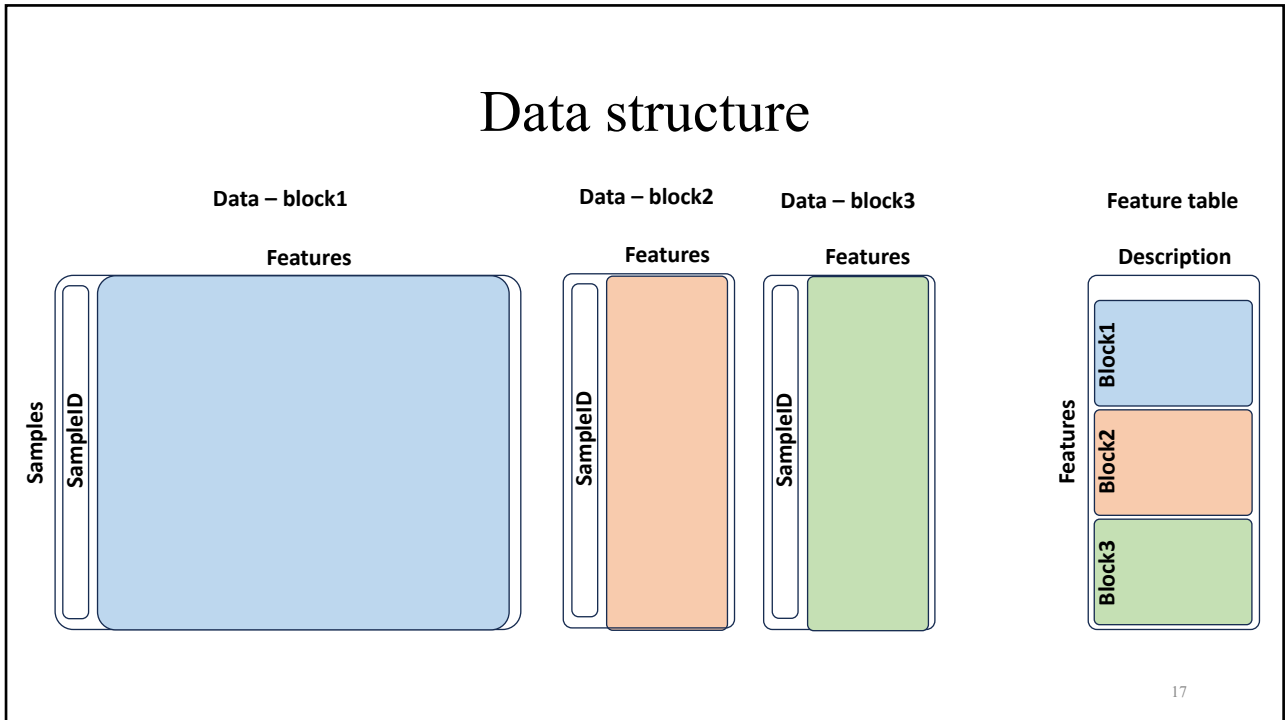
Features

Description

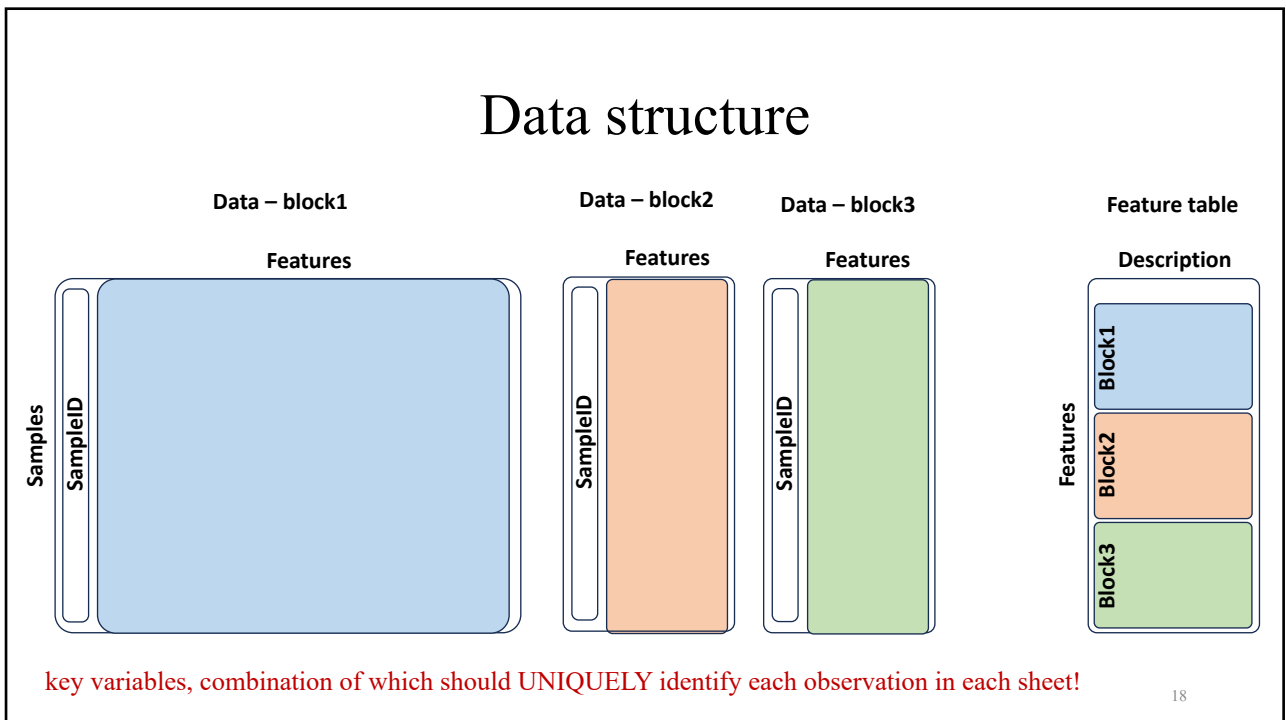


16

16



17



18

Wide format of dependent data

e.g. different timepoints for the same variable

FD	FE	FF	FG
ADL M6	IADL M6	ADL6 M12	IADL4 M12
3	0	3	0

19

19

Wide format of dependent data

e.g. different timepoints for the same variable

FD	FE	FF	FG
ADL M6	IADL M6	ADL6 M12	IADL4 M12
3	0	3	0

20

20

Long format



	A	B	FQ	FR	FS	FT
1	PatientCode	Timepoint	nCD4	%CD4	nCD8	%CD8
2	HIP001	J0	498,64	51,2	311,31	31,96
3	HIP001	J10	666,36	48,59	471,77	34,4
4	HIP001	J13	810,69	50,11	538,75	33,3
5	HIP001	J3	675,57	47,07	530,55	36,96
6	HIP001	J4	419,65	43,48	362,09	37,52
7	HIP001	J7	766,91	47,76	555,47	34,6
8	HIP001	M17	730,35	49,65	482,35	32,79
9	HIP002	J0	251,25	54,8	123,56	26,95
10	HIP002	J1	614,52	62,24	278,86	28,25
11	HIP002	J3	730,64	58,09	348,43	27,7
12	HIP002	J6	516,98	64,84	219,41	27,52
13	HIP002	M15	633,34	60,27	272,96	25,97

21

21

Other remarks

- Missing values are marked as ` `', `?', `ND`, `nd`, `NA`, `na` or `#N/A` (excel-integrated NA notation)
- Check attentively all missing values for response variables (they will be ignored by the model if you don't manage it yourself)
- Models require variability inside each group

22

22

HipAge1_original

PatientCode	Timepoint	Follow up arrival	Group	Group bis	StudyDSIncl	DCD	Age	sexe 1-h	delai hop-g	delai hop c	duree chir	Corticoides	Données orthopédie	Fracture	Type fractu	Petrochan	Cervicale	Fract patho	Polytraum	COMORBIDITES	Demence	MMS ant	Depression	MDPI	HTA	Diab
HIP001	J0	inter	HIP	HIP	TRUE	0	90	0	24	54	120	0	garden 3 gs	Garden 4	0	1	0	0	0	0	1	22	0	0	0	
HIP001	J10	inter	HIP	HIP	TRUE	0	90	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP001	J13	sortie	HIP	HIP	TRUE	0	90	0	0	24	54	120	0	garden 3 gs	Garden 4	0	1	0	0	0	1	22	0	0	0	
HIP001	J4	post	HIP	HIP	TRUE	0	90	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP001	J7	inter	HIP	HIP	TRUE	0	90	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP001	M17	Ig terme	HIP	HIP	TRUE	0	91	0	0	18	13	130	0	petroch	Petrochan	1	0	0	0	0	1	19	0	0	1	
HIP002	J0	pre	HIP	HIP	TRUE	0	91	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP002	J1	post	HIP	HIP	TRUE	0	91	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP002	J3	inter	HIP	HIP	TRUE	0	91	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP002	J6	sortie	HIP	HIP	TRUE	0	91	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP002	M15	Ig terme	HIP	HIP	TRUE	0	92	0	0	0	0	0	0	0	0	0	0	0	0	0	1	24	1	0	1	
HIP003	J0	X	C	CTRL	TRUE	0	80	0	0	0	0	0	0	Cognitif	HdJ	0	0	0	0	0	0	16	0	0	1	
HIP004	J0	pre	HIP	HIP	TRUE	1	82	0	25	17	300	1	fracture p	Garden 4	0	1	1	1	1	0	0	1	0	0	0	
HIP004	J6	post	HIP	HIP	TRUE	1	82	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP004	J6	sortie	HIP	HIP	TRUE	1	82	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP005	J0	pre	HIP	HIP	TRUE	0	88	0	22	39	120	0	garden 1 vi	Garden 1	0	1	0	0	0	1	15	0	0	0	1	
HIP005	J1	pre	HIP	HIP	FALSE	0	88	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP005	J1	post	HIP	HIP	TRUE	0	88	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP005	J2	inter	HIP	HIP	FALSE	0	88	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP005	J5	sortie	HIP	HIP	TRUE	0	88	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP006	J0	X	C	CTRL	FALSE	0	86	0	0	0	0	0	0	Cognitif	HdJ	0	0	0	0	0	0	16	0	0	1	
HIP007	J0	pre	HIP	HIP	TRUE	0	93	0	18	5	110	0	garden 1 vi	Garden 1	0	1	0	0	0	1	18	0	0	0	1	
HIP007	J5	post	HIP	HIP	TRUE	0	93	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP007	J7	sortie	HIP	HIP	TRUE	0	93	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP008	J0	X	C	CTRL	TRUE	0	99	1	0	0	0	0	0	Dyspnée fai	HdJ	0	0	0	0	0	25	0	0	0	0	
HIP009	J0	X	C	CTRL	TRUE	0	86	1	0	0	0	0	0	Cognitif	HdJ	0	0	0	0	0	26	1	0	1	1	
HIP010	J0	X	C	CTRL	FALSE	0	90	1	0	0	0	0	0	Dysgnie	HdJ	0	0	0	0	0	1	0	0	0	0	
HIP011	J0	X	C	PB	TRUE	0	81	1	0	0	0	0	0	Chute avec	UGA	0	1	0	0	0	0	0	0	0	1	
HIP012	J0	X	C	CTRL	TRUE	0	89	0	0	0	0	0	0	Bilan Insuff	HdJ	0	0	0	0	1	18	0	0	1	1	
HIP013	J0	X	C	CTRL	TRUE	0	82	0	0	0	0	0	0	Bilan AEG	HdJ	0	0	0	0	0	29	1	0	1	1	
HIP014	J0	X	C	CTRL	TRUE	0	79	1	0	0	0	0	0	Cognitif	HdJ	1	1	1	1	1	25	1	0	0	0	
HIP015	J0	X	C	CTRL	TRUE	0	85	1	0	0	0	0	0	Cognitif	HdJ	1	1	1	1	1	17	1	1	1	1	
HIP016	J0	X	C	CTRL	TRUE	0	90	0	0	0	0	0	0	Cognitif	HdJ	1	1	1	1	1	23	1	0	0	0	
HIP016	M3	X	C	CTRL	FALSE	0	91	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP017	J0	pre	HIP	HIP	TRUE	1	83	1	48	24	160	1	0	Garden 4	0	1	0	0	0	0	26	1	0	1	1	
HIP017	J6	post	HIP	HIP	TRUE	1	83	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP017	J9	sortie	HIP	HIP	TRUE	1	83	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP018	J0	pre	HIP	HIP	TRUE	0	105	0	24	9	110	1	0	Petrochan	1	0	0	0	0	0	0	23	0	0	1	
HIP018	J2	post	HIP	HIP	TRUE	0	105	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP018	M13	Ig terme	HIP	HIP	TRUE	0	106	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP018	M6	Ig terme	HIP	HIP	TRUE	0	106	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
HIP019	J0	X	C	CTRL	TRUE	0	83	0	0	0	0	0	0	Cognitif	HdJ	0	0	0	0	0	1	26	1	0	0	

23

HipAge1_original

ColNames	Categories	MeanBol	Ignore	IgnoreLate	Known_Var	Chi2
PatientCode	ID	NA	0	0	0	FALSE
Timepoint	ID	NA	0	0	0	FALSE
Follow up	ID	NA	0	0	0	FALSE
Group	ID	NA	0	0	0	FALSE
Group bis	ID	NA	0	0	0	FALSE
DCD	ID	NA	0	0	0	FALSE
Age	Clinic	1	0	0	0	FALSE
sexe 1-h	Clinic	0	0	0	0	TRUE
delai hop-geriatrie (h)	Clinic	0	1	1	1	FALSE
delai hop chirurgie (h)	Clinic	0	1	1	1	FALSE
duree chirurgie (min)	Clinic	0	1	1	1	FALSE
Corticoides au bloc	Clinic	0	1	1	1	TRUE
Données orthopédie	Clinic	0	1	1	1	TRUE
Fracture	Clinic	0	1	1	1	TRUE
Type fracture	Clinic	0	1	1	1	FALSE
Petrochanter	Clinic	0	1	1	1	FALSE
Cervicale	Clinic	0	1	1	1	FALSE
Fract pathologique	Clinic	0	1	1	1	FALSE
Polytraum	Clinic	0	1	1	1	FALSE
COMORBIDITES	Clinic	0	1	1	1	FALSE
Demence	Clinic	1	0	0	0	TRUE
MMS ant	Clinic	0	0	0	0	FALSE
Depression	Clinic	1	0	0	0	TRUE
MDPI	Clinic	1	1	1	1	FALSE
HTA	Clinic	1	1	1	1	FALSE
Diabete	Clinic	1	1	1	1	TRUE
ACFA	Clinic	1	1	1	1	TRUE
Obesite	Clinic	1	1	1	1	TRUE
Ins Coro	Clinic	1	1	1	1	TRUE
Ins card	Clinic	1	1	1	1	TRUE

24

HipAge1_Filtered_and_fixed

PatientCode	DCD	SurvivalTime	Age	sexe 1=h	Demence	MMS ant	cancer	ASA	CIRS 52	Charlson pi	Rockwood	IADL 4 ant	ADL6 ant	Hb prog	total culoti	albumine	Poids	Creatinine	Cockroft	vitamine D	folate	B12	HSV titres	VZV titres	CMV titres	CRP
HIP001	0	365	90	0	1	22	1	2	14	10	6	1	4.5	13	1	26	53	51	65	10	45	388	3.73	3.52	236	
HIP002	0	365	91	0	1	19	0	3	15	10	7	1	5.5	12.8	0	32	62.8	130	25	10	45	388	3.73	2.31	236	
HIP004	1	197	82	0	0	1	1	2	7	6	3	4	6	11.6	0	26	50	24	126	24	13	952	2.29	1.94	1.07	
HIP005	0	729	88	0	1	15	0	3	8	7	3	2	6	11	0	34	52	82	41	15	10	327	2.34	1.94	1.07	
HIP007	0	365	93	0	1	18	0	3	10	8	7	1	3	12	0	27	42	68	43	15	10	334	1.94	1.07	1.07	
HIP017	1	135	83	1	0	26	0	3	14	7	5	3	6	13.1	0	33	76	136	39	18	227	4.01	1.94	1.07		
HIP018	0	365	105	0	0	23	0	3	11	10	4	1	3	14.4	2	27	43	81	20	15	10	342	1.07	1.07	1.07	
HIP021	1	18	93	0	0	1	1	2	5	7	5	1	3.5	13.6	0	23	45	46	45	22	14	294	3.73	3.18	1.07	
HIP022	0	365	92	1	0	1	1	2	14	9	7	1	1	14.9	1	33	55	75	54	47	15	10	3.73	0.91	1.07	
HIP023	0	365	98	0	1	23	1	3	21	12	7	0	3.5	13.2	2	30	68	67	44	15	10	4.06	2.22	1.07		
HIP024	0	30	95	1	0	1	1	2	7	8	6	0	0.5	14	1	28	45	119	23	15	10	1.94	1.07	1.07		
HIP025	0	365	80	0	1	1	1	2	9	5	4	4	6	13.4	0	39	59.3	48	77	15	10	1.94	1.07	1.07		
HIP038	1	365	93	0	1	24	1	3	14	13	6	1	5	12.6	1	31	49	57	53	43	8.79	264	1.94	1.07	1.07	
HIP036	1	37	95	1	0	1	1	2	11	10	3	2	5	13.1	0	33	92.5	107	47	62	14.8	252.5	2.76	1.97	1.07	
HIP038	0	365	87	1	0	1	1	2	8	3	3	4	6	13.5	0	34	64	43	103	15	10	2.34	1.97	1.07		
HIP042	0	365	74	0	1	1	1	2	8	3	3	4	6	13.5	0	34	64	43	103	15	10	2.34	1.97	1.07		
HIP044	0	365	95	0	1	1	1	2	6	8	6	0	2.5	13.2	0	34	45	56	37	15	10	1.94	1.07	1.07		
HIP052	0	365	80	1	0	1	1	2	6	8	6	0	2.5	13.2	0	34	45	56	37	15	10	1.94	1.07	1.07		
HIP057	0	365	75	0	1	1	1	2	6	4	1	4	6	13.5	0	37	57	50	77	37	37	677	4.55	1.94	1.07	
HIP064	0	365	80	0	1	1	1	2	6	5	3	3	5.5	13.9	0	36	50	58	55	15	302	-0.05	1.94	1.07		
HIP067	0	1152	86	0	1	1	1	2	13	9	5	4	5.5	10.2	1	24	55	75	42	27	15	10	34	1.94	1.07	1.07
HIP068	1	34	87	0	0	24	0	3	19	8	4	4	5.5	12.3	1	24	27	62	24	32	15	10	3.64	5.43	1.07	
HIP073	0	365	78	0	1	1	1	2	8	9	5	3	5.5	13.8	0	32	67	67	48	26	15	10	3.98	2.3	68	
HIP074	0	365	85	1	0	1	1	2	13	5	6	0	3.5	12	0	34	84	81	70	15	10	1.94	1.07	1.07		
HIP075	0	365	91	0	1	1	1	2	7	6	6	0	6	11.4	1	31	49	57	53	43	8.79	264	1.94	1.07	1.07	
HIP083	1	80	78	1	0	22	1	3	26	15	5	1	2.5	12.8	0	29	62	160	28	15	10	1.94	1.07	1.07		
HIP084	0	365	89	0	1	26	1	3	8	9	3	4	6	14.8	0	27	53	56	53	18	450	1.94	1.07	1.07		
HIP091	0	803	88	1	0	1	1	2	8	7	5	0	3	13.8	0	27	79	95	54	15	10	1.94	1.07	1.07		
HIP093	0	365	100	0	1	1	1	2	10	9	4	4	5.5	11.7	0	26	56	83	30	28	15	10	6.33	2.52	1.07	
HIP098	0	365	85	0	1	1	1	2	8	7	2	4	5.5	11.5	2	31	58	76	43	15	10	1.94	1.07	1.07		
HIP109	0	365	76	0	1	1	1	2	10	5	6	0	5	14.7	0	25	56	60	60	15	10	5.14	3.54	2.27	1.07	
HIP110	0	365	91	0	1	23	0	3	6	7	6	3	4.5	12.2	2	21	52	56	46	11.7	151	1.94	1.07	1.07		
HIP113	0	365	85	1	0	29	0	3	6	8	4	4	5.5	11.7	0	27	88	75	53	15	10	6.4	3.45	1.07		
HIP116	0	365	91	0	1	1	1	2	11	7	6	3	5.5	13.7	0	28	51	46	55	15	10	1.94	1.07	1.07		
HIP119	0	365	84	0	1	1	1	2	8	5	3	4	6	11.9	0	30	63	64	62	25	15	10	5.79	1.57	1.07	
HIP120	0	365	91	0	1	23	0	3	14	11	6	1	3.5	11.8	0	26	57	67	48	15	10	1.94	1.07	1.07		
HIP121	0	365	89	0	1	1	1	2	8	4	4	4	5.5	14.6	0	29	93	64	75	8	13	322	1.94	1.07	1.07	
HIP122	0	365	79	0	1	23	0	3	8	4	5	2	5	13.1	0	25	79	72	68	24	6.6	319	5.686	1.94	1.07	
HIP123	0	365	84	0	1	1	1	2	10	11	6	2	5	9.7	2	32	45	73	34	0	4.8	308	-0.24	2.19	250	
HIP124	0	365	85	1	0	1	1	2	7	5	6	1	4	11.3	6	23	57	54	71	123	8.9	-0.1	2.52	1.07	1.07	
HIP125	0	365	90	1	0	1	1	2	2	0	6	3	5.5	15.6	1	24	91	42	133	21	9.7	701	1.94	1.07	1.07	
HIP126	1	19	96	0	1	10	1	3	13	9	5	0	3.5	11	2	23	33	30	66	24	25.9	168	1.94	1.07	1.07	
HIP128	0	365	83	0	1	1	1	2	9	8	3	4	6	11.2	1	25	66	85	47	14	15	276	1.94	1.07	1.07	

25

rshiny.immublab.fr

Feature Selector

Download outputs: select all deselect all

Choose excel file: HipAge1_Filtered_and_fixed.xlsx

Sheet(s) with the source data: Source Data

Sheet with the information about variables: Variable Information

Select the "variable name" column: ColNames

Submit

Overview of the source data

PatientCode	DCD	SurvivalTime	Age	sexe 1=h	Demence	MMS ant	cancer	ASA	CIRS 52	Charlson pondt
1 HIP001	0	365	90	0	1	22	1	2	14	
2 HIP002	0	365	91	0	1	19	0	3	15	
3 HIP004	1	197	82	0	0	1	1	2	7	
4 HIP005	0	729	88	0	1	15	0	2	8	
5 HIP007	0	365	93	0	1	18	0	3	10	
6 HIP017	1	135	83	1	0	26	0	2	14	
7 HIP018	0	365	105	0	0	23	0	3	11	
8 HIP021	1	18	93	0	0	0	0	2	5	
9 HIP022	0	365	92	1	0	0	0	3	14	
10 HIP023	0	365	98	0	1	23	1	3	21	
11 HIP024	0	30	95	1	1	0	0	3	7	

Showing 1 to 60 of 60 entries

26

26

Feature Selector

Download outputs:
 select all
 deselect all

Choose excel file
 Browse... HipAge1_Filtered_and_fixed.xlsx
 Upload complete

Sheet(s) with the source data:
 Source Data

Sheet with the information about variables:
 Variable Information

Select the "variable name" column:
 ColNames

Submit

Overview of the source data Overview of the information about variables

Search:

	ColNames	Categories	MeanBoI	Ignore	IgnoreLateMarker	Known_Var	Chi2
1	PatientCode	ID		0	0	0	false
2	Timepoint	ID		0	0	0	false
3	Follow up	ID		0	0	0	false
4	Group	ID		0	0	0	false
5	Group bis	ID		0	0	0	false
6	DCD	ID		0	0	0	false
7	Age	Clinic	1	0	0	0	false
8	sexe 1=h	Clinic	0	0	0	1	true
9	delai hop-geriatrie (h)	Clinic	0	1	1	1	false
10	delai hop chirurgie (h)	Clinic	0	1	1	1	false
11	duree chirurgie (min)	Clinic	0	1	1	1	false

Showing 1 to 277 of 277 entries

27

27

Clean the data

- Removes variables with no variation
- Removes dates
- Sets appropriate types to variables (categorical / numeric / logical)
- Problem – if looks numeric, but is to be better considered categorical
 e.g. medical score going from 1 to 3 is either approximated as a number or as the group variable (giving 3 groups of patients). For small numbers of groups and in the absence of numerical relations between the score values.
- The user has to choose two heuristics – maximal nb of categories and digits for such variables to be considered as categorical

28

28

For a numeric looking variable to be considered as categorical

Maximal number (included) of levels:

Maximal number (included) of digits:

Automatic changes to the data
Read it attentively and modify heuristics to obtain the data that makes sense

Save

Search:

	Variable	Change
1	PatientCode	is set to factor
2	DCD	is set to factor
3	SurvivalTime	is left as numeric
4	Age	is left as numeric
5	sexe 1=h	is set to factor
6	Demence	is set to factor
7	MMS ant	is left as numeric
8	cancer	is set to factor
9	ASA	is set to factor
10	CIRS 52	is left as numeric
11	Charlson pondéré	is left as numeric
12	Rockwood	is left as numeric
13	IADL 4 ant	is left as numeric
14	ADL6 ant	is left as numeric
15	Hb preop	is left as numeric

29

Filter

- Select the ID variable (uniquely identifying each filtered observation)
- Select the response variable(s)
- Set filters for response variable(s) and predictors (same or separate)
- If separate filters are set, pay attention to whether you consider the response variable as the potential predictor or not ! (e.g. we might be interested to predict post-surgery CRP level with pre-surgery CRP level, but not with post-surgery CRP)

30

30

Feature Selector

Select the ID variable: PatientCode

Select all response variables: DCD SurvivalTime

Filter rows (observations) for response variables

Filter 1: Include

Filter 2: AND

Follow up: ppi

Remove Filter Add Filter

for predictors: Apply same filters

Filter columns (variables): Filter 1: Include

IgnoreLateMarker: 0

Filter 2: AND

Exclude: Categories: ID

Remove Filter Add Filter

Select variables to exclude from the analysis (but submit filters first): Denutrition Dindo-Clavien

Select variables to include anyway (but submit filters first):

Submit All

Overview of the selected data Visualize NA's

Save

PatientCode	DCD	SurvivalTime	Age	sexe 1=h	Demence	MMS ant	cancer	ASA	CIRS 52	Charlson pondri	Rockwood	IADL 4= ant	ADL6 ant	Hb prop	total culets	albumine	Polds	Creatinine	
1 HIP001	0		90	0	1	22	1	2	14	10		6	1	4.5	13	1	26	53	1
2 HIP002	0		91	0	1	19	0	3	15	10		7	1	5.5	12.8	0	32	62.8	11
3 HIP004	1	197	82	0	0		1	2	7	6		3	4	6	11.6	0	26	50	1
4 HIP005	0	729	88	0	1	15	0	2	8	7		3	2	6	11	0	34	52	1
5 HIP007	0		93	0	1	18	0	3	10	8		7	1	3	12	0	27	42	1
6 HIP017	1	135	83	1	0	26	0	2	14	7		5	3	6	13.1	0	33	76	11
7 HIP018	0	105	0	0	0	23	0	3	11	10		4	1	3	14.4	2	27	43	1
8 HIP021	1	18	93	0	0	0	2	5	7	5		5	1	3.5	13.6	0	23	45	1
9 HIP022	0	92	1	0	0	0	3	14	8	7		7	1	1	14.9	1	33	55	1
10 HIP023	0		98	0	1	23	1	3	21	12		7	0	3.5	13.2	2	30	68	1
11 HIP024	0	30	95	1	1	0	3	7	8	6		6	0	0.5	14	1	28	45	11
12 HIP025	0		80	0	0	0	2	9	5	4		4	4	6	13.4	0	39	59.3	1
13 HIP033	0		93	0	1	24	1	3	14	13		6	1	5	12.6	1	31	49	1
14 HIP036	1	37	95	1	0	0	4	11	10	10		3	2	5	13.1	0	33	92.5	11

Showing 1 to 60 of 60 entries

NA's in response variables: SurvivalTime (Check the selection of the source data, otherwise those observations will be ignored.)

31

Feature Selector

Select the ID variable: PatientCode

Select all response variables: DCD SurvivalTime

Filter rows (observations) for response variables

Filter 1: Include

Remove Filter Add Filter

for predictors: Apply same filters

Filter columns (variables):

Select variables to exclude from the analysis (but submit filters first):

Select variables to include anyway (but submit filters first):

Submit All

Overview of the selected data Visualize NA's

Save

PatientCode	DCD	SurvivalTime	Age	sexe 1=h	Demence	MMS ant	cancer	ASA	CIRS 52	Charls pondri
1 HIP001	0	365	90	0	1	22	1	2	14	
2 HIP002	0	365	91	0	1	19	0	3	15	
3 HIP004	1	197	82	0	0		1	2	7	
4 HIP005	0	729	88	0	1	15	0	2	8	
5 HIP007	0	365	93	0	1	18	0	3	10	
6 HIP017	1	135	83	1	0	26	0	2	14	
7 HIP018	0	365	105	0	0	23	0	3	11	
8 HIP021	1	18	93	0	0		0	2	5	
9 HIP022	0	365	92	1	0		0	3	14	
10 HIP023	0	365	98	0	1	23	1	3	21	
11 HIP024	0	30	95	1	1		0	3	7	
12 HIP025	0	365	80	0	0		0	2	9	
13 HIP033	0	365	93	0	1	24	1	3	32	14
14 HIP036	1	37	95	1	0		0	4	11	

32

33

Missing Values

Stringent = exclude rows / columns having too many NA's

Impute (= fill in NA's in each variable) with :

- central tendency – the mean, median or mode
- k Nearest Neighbors – the mean of k most resembling samples
- other methods are to come...

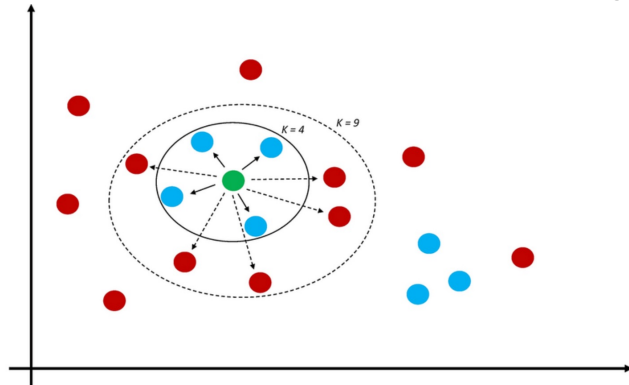
Musolf, Anthony & Holzinger, Emily & Malley, James & Bailey-Wilson, Joan. (2022). What makes a good prediction? Feature importance and beginning to open the black box of machine learning in genetics. Human Genetics. 141. 10.1007/s00439-021-02402-z.

34

34

Missing Values

kNN – replace NA's with the mean of k most resembling samples



Musolf, Anthony & Holzinger, Emily & Malley, James & Bailey-Wilson, Joan. (2022). What makes a good prediction? Feature importance and beginning to open the black box of machine learning in genetics. Human Genetics. ³⁵

35

36

Stringent

Maximal tolerated % of NA by row: 60.8

Maximal tolerated % of NA by variable (column): 60.6

Imputation

Select imputation method: k Nearest Neighbors

Choose number of neighbors: 5

Save

png jpg svg pdf

Observations

Missing (10.9%) Present (89.1%)

Variable	% NA
1 VZV titre	54
2 vitamine D (ng/ml)	46
3 HSV titre	37

37

Imputation effect

before after

value

1.0
0.5
0.0
-0.5

Save imputation histograms (it may take a lot of time to generate)

38

The screenshot shows the 'Feature Selector' software interface. On the left, there are sliders for 'Maximal tolerated % of NA by row:' (set to 60.8) and 'Maximal tolerated % of NA by variable (column):' (set to 60.6). Below these are imputation settings: 'Select imputation method:' set to 'k Nearest Neighbours' and 'Choose number of neighbors:' set to 5. On the right, there are tabs for 'Stringent', 'Worst variables after stringent', 'Imputation effect on correlations', and 'Overview of Imputed Data'. A table of 14 patient records is displayed with columns: PatientCode, DCD, SurvivalTime, Age, sexe 1=h, Demence, cancer, ASA, CIRS 52, Charlson pondéré, and Rc.

PatientCode	DCD	SurvivalTime	Age	sexe 1=h	Demence	cancer	ASA	CIRS 52	Charlson pondéré	Rc
1	HIP001	0	365	90	0	1	2	14	10	
2	HIP002	0	365	91	0	1	3	15	10	
3	HIP004	1	197	82	0	0	2	7	6	
4	HIP005	0	729	88	0	1	2	8	7	
5	HIP007	0	365	93	0	1	3	10	8	
6	HIP017	1	135	83	1	0	2	14	7	
7	HIP018	0	365	105	0	0	3	11	10	
8	HIP021	1	18	93	0	0	2	5	7	
9	HIP022	0	365	92	1	0	3	14	9	
10	HIP023	0	365	98	0	1	3	21	12	
11	HIP024	0	30	95	1	1	3	7	8	
12	HIP025	0	365	80	0	0	2	9	5	
13	HIP033	0	365	93	0	1	3	14	13	
14	HIP036	1	37	95	1	0	4	11	10	39

39

Odds Ratios

Preliminary analysis in case of binary response

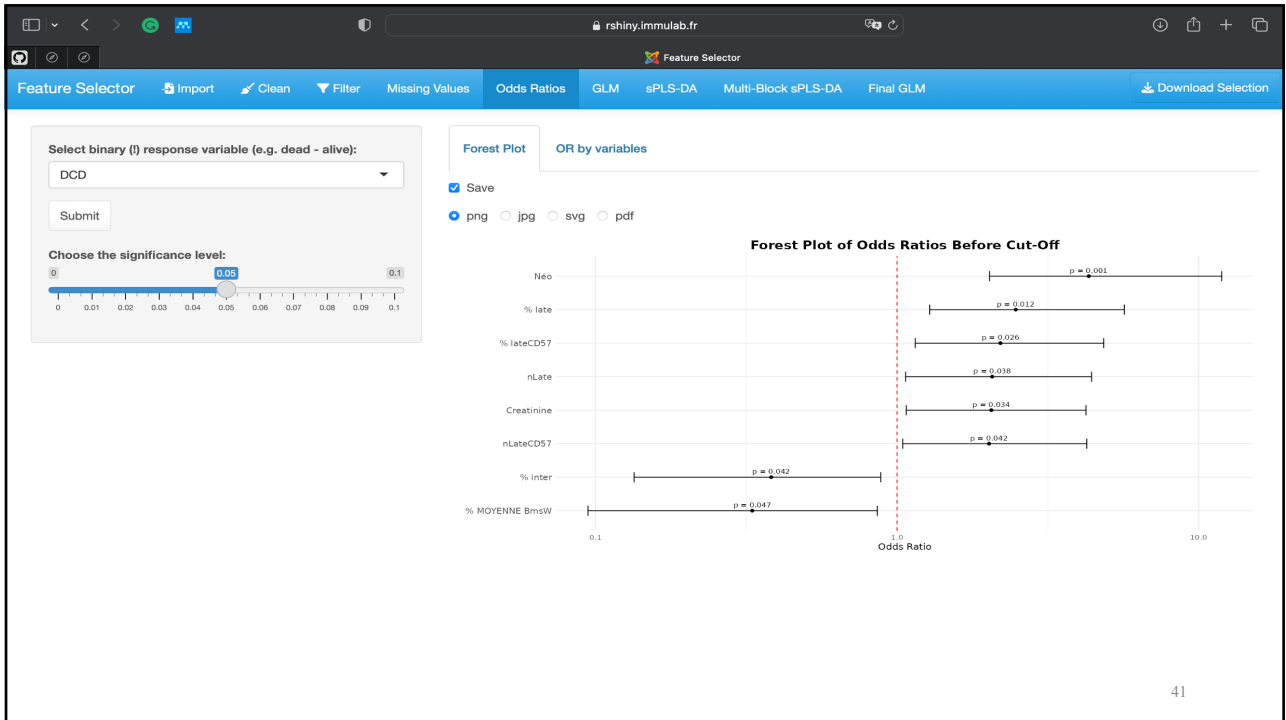
Variables are considered one by one, interactions are not taken into account

Forest Plot of Odds Ratios Before Cut-Off

The plot shows Odds Ratios on a logarithmic x-axis (0.1 to 10.0) with a vertical dashed red line at 1.0. Variables and their p-values are: Néo (p=0.001), % late (p=0.012), % lateCD57 (p=0.026), nLate (p=0.038), Creatinine (p=0.034), nLateCD57 (p=0.042), % inter (p=0.042), and % MOYENNE BmsW (p=0.047).

attention to the x axis!

40



41

Search:

predictor	OR	lower	upper	p-value
1 Age	1.1517	0.6042	3.8229	0.7507
2 sexe 1=h	3.8	0.9001	16.375	0.0657
3 Demence	0.25	0.0352	1.1239	0.0995
4 MMS ant	0.7643	0.2289	2.5996	0.6423
5 cancer	1.6286	0.3091	7.1093	0.5298
6 ASA	8508962.397	0		0.9936
7 CIRS 52	1.8836	0.9975	3.8769	0.0574
8 Charlson pondéré	1.9003	0.9368	4.2356	0.0878
9 Rockwood	0.8117	0.398	1.6788	0.5609
10 IADL 4 ant	0.757	0.3454	1.5823	0.4583
11 ADL6 ant	0.7299	0.3756	1.4758	0.3535
12 Hb preop	0.6616	0.3264	1.3274	0.233
13 total culots	0.9507	0.4003	1.8754	0.8937
14 albumine	0.502	0.1946	1.1003	0.1112
15 Poids	1.1763	0.5827	2.3557	0.6491

42

GLM

Multiple Regression models can predict :

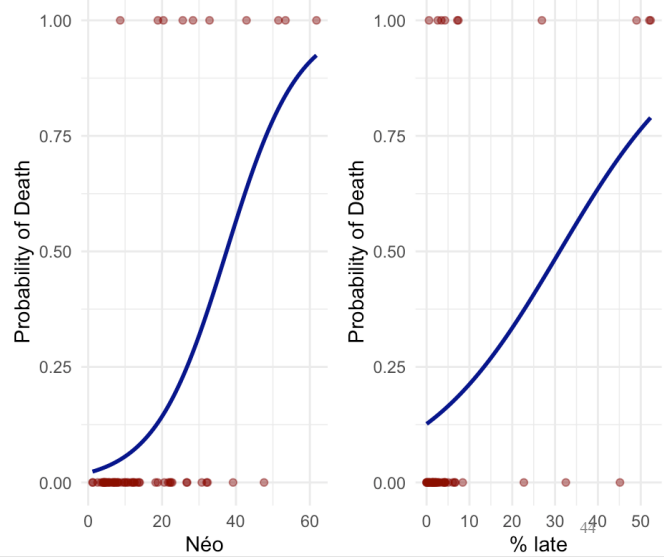
- Linear – numeric outcome
- Logistic – binary outcome
- Cox Proportional Hazards model (Survival analysis) – hazard function

43

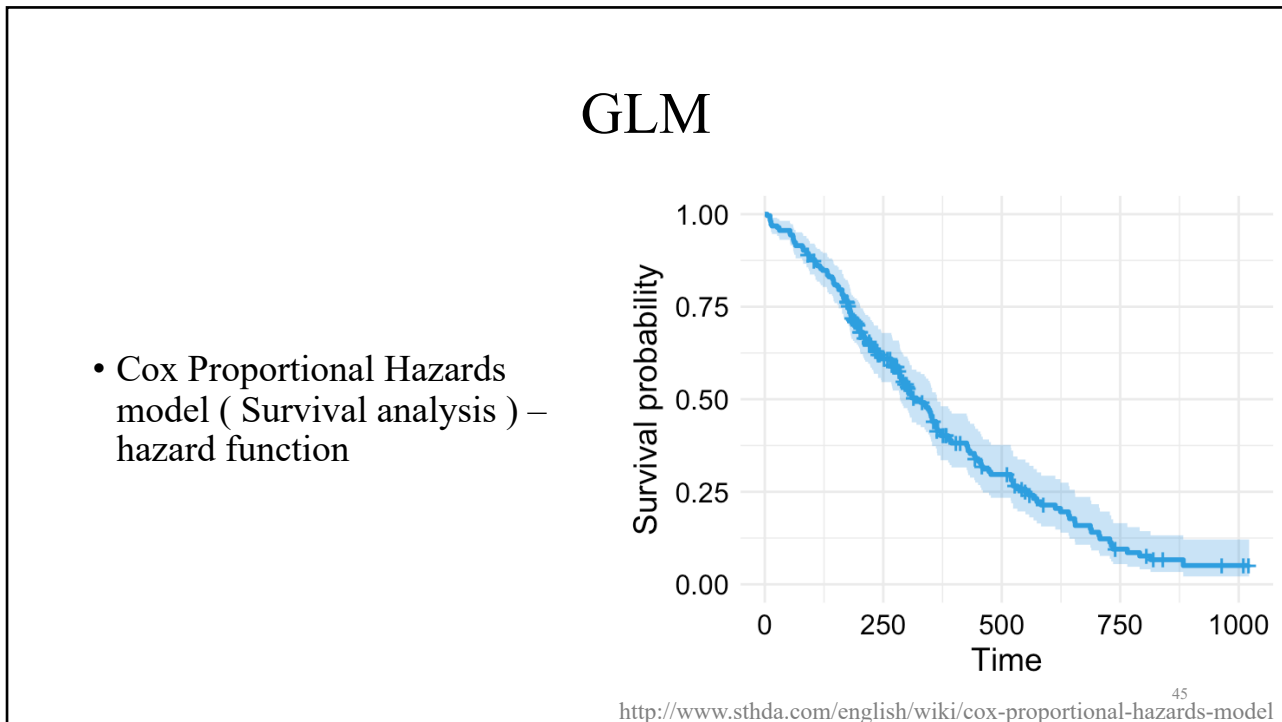
43

GLM

- Logistic – binary outcome



44



45

GLM Assumptions

Linear regression:

- **linearity** of relationship
- independence of observations
- **no perfect multicollinearity**
- (no endogeneity)

Regularization is more robust to violations of these assumptions:

- homoscedasticity (constant variance)
- residual normality

46

46

GLM Assumptions

Logistic regression:

- **linearity** of log-odds
- independence of observations
- **no perfect multicollinearity**

- large sample size

47

47

GLM Assumptions

Cox PH model:

- **linearity** of log-hazard
- independence of observations
- **no perfect multicollinearity**

- proportional hazards (hazard ratios are constant over time)
- non-informative censoring

48

48

Regularized GLM

Cost Function = Loss Function + Penalty

Lasso (L¹ Norm)

$$\lambda \sum_{j=1}^p |\beta_j|$$

ElasticNet

$$\lambda \left[\alpha \sum_{j=1}^p |\beta_j| + \frac{(1-\alpha)}{2} \sum_{j=1}^p \beta_j^2 \right]$$

Ridge (L² Norm)

$$\frac{\lambda}{2} \sum_{j=1}^p \beta_j^2$$

Impose alpha

Choose the alpha (0=Ridge, 1=Lasso) to impose:

0 1

49

49

Regularized GLM Assumptions

- **the same scale** – variables must be standardized, but no clear consensus on whether dummy variables (indicators of levels of categorical variables) have to be standardized or not – try both ways

Standardize dummies

50

50

Model Robustness

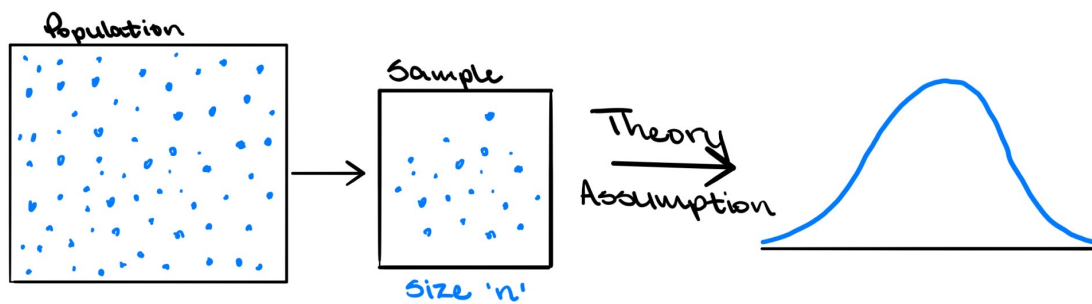
If I repeat the experiment, will I get the same result?

51

51

Model Robustness

Theoretical approach:

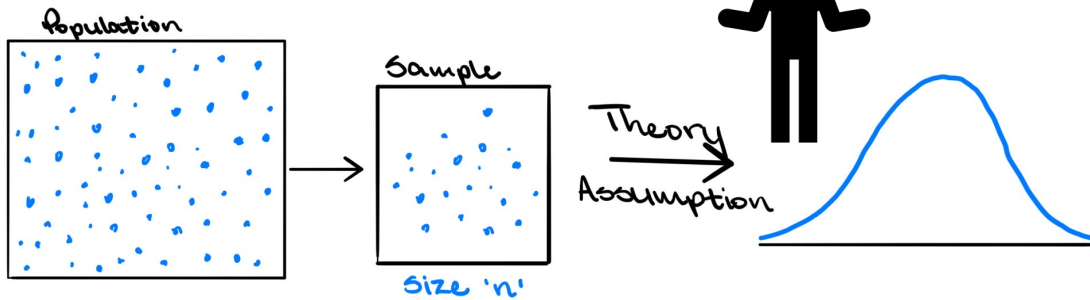


<https://www.linkedin.com/pulse/bootstrapping-statistics-what-why-its-used-trist-n-joseph>⁵²

52

Model Robustness

Theoretical approach:

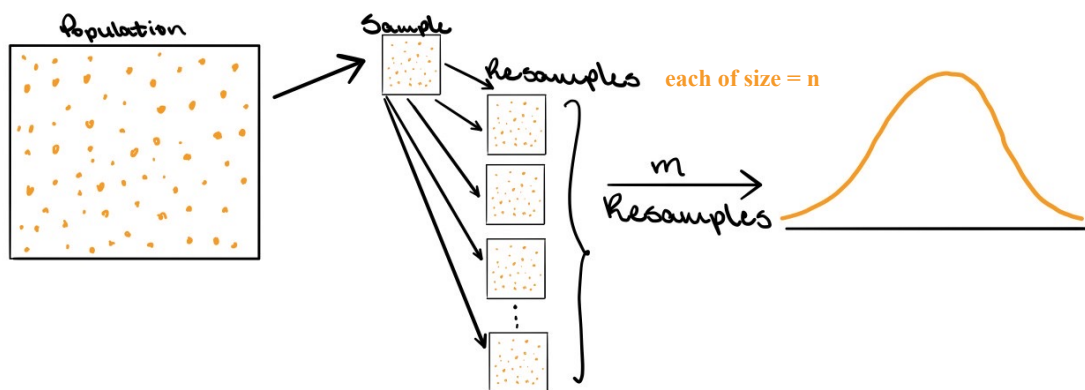


<https://www.linkedin.com/pulse/bootstrapping-statistics-what-why-its-used-trist-n-joseph>⁵³

53

Model Robustness

Brute force simulation approach: **Bootstrapping** or Cross-Validation



<https://www.linkedin.com/pulse/bootstrapping-statistics-what-why-its-used-trist-n-joseph>⁵⁴

54

Feature Selector

Select the type of the model:
 Logistic Regression
 Linear Regression
 Logistic Regression
 Cox Proportionnal Hazards

Impose alpha

Choose the alpha (0=Ridge, 1=Lasso) to impose:
 0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

Standardize dummies

Choose the number of bootstraps:
 1 100 1,000

Choose the oversampling:
 1 1.4 1.8 2.2 2.6 3 3.4 3.8 4.2 4.6 5

Submit

Choose the threshold bootstrap frequency for the plot:
 0 0.1 1

Bootsrapped plot | Bootsrapped FS table | FS on the original data

Save

png jpg svg pdf

55

55

Feature Selector

Select the type of the model:
 Cox Proportionnal Hazards

Select the name of the binary status response variable:
 DCD

Select the name of the survival time variable:
 SurvivalTime

Impose alpha

Choose the alpha (0=Ridge, 1=Lasso) to impose:
 0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

Standardize dummies

Choose the number of bootstraps:
 1 100 1,000

Choose the oversampling:
 1 1.4 1.8 2.2 2.6 3 3.4 3.8 4.2 4.6 5

Submit

Choose the threshold bootstrap frequency for the plot:
 0 0.1 1

Bootsrapped plot | Bootsrapped FS table | FS on the original data

Save

png jpg svg pdf

Selected features coming up in more than 10% of bootstrapped models

Frequency

Feature	Frequency
Néo	0.85
% MOYENNE BmsW	0.75
% late	0.65
Creatinine	0.60
Poids	0.55
%CD8M CD57+	0.50
% inter	0.45
Eosino	0.40
IL1B	0.35
nLate	0.30
%CD8M Ki67+	0.25
vitamine D (ng/ml)	0.20
nNKCD3+	0.15
IP10	0.10

56

56

The screenshot shows the 'Feature Selector' web application interface. On the left, the configuration panel is set to 'Cox Proportionnal Hazards' model with response variable 'DCD' and survival time variable 'SurvivalTime'. The 'Impose alpha' checkbox is checked, and the alpha slider is set to 1. The number of bootstraps is 100, and the oversampling factor is 2. The threshold bootstrap frequency for the plot is 0.1. On the right, the 'Bootsrapped FS table' is displayed, showing a list of features and their frequencies.

Feature	Frequency	
80	Néo	0.88
63	% MOYENNE BmsW	0.38
17	Creatinine	0.31
58	% late	0.31
16	Poids	0.3
54	%CD8M CD57+	0.29
57	% inter	0.28
42	Eosino	0.27
73	IL1B	0.26
94	nLate	0.22
55	%CD8M Ki67+	0.18
19	vitamine D (hg/ml)	0.14
35	nNKCD3+	0.12
77	IP10	0.11
66	% CD14 +	57 0.09

57

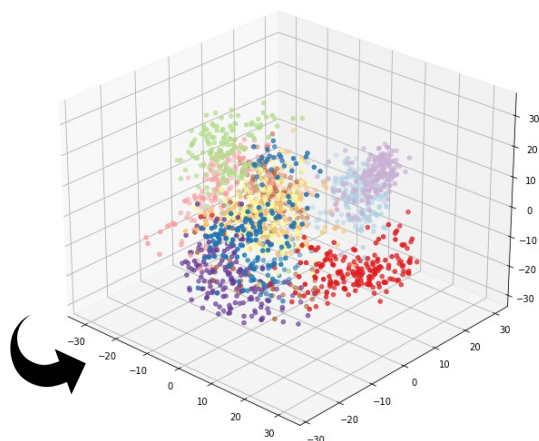
This screenshot shows the same 'Feature Selector' interface as above, but with the 'FS on the original data' tab selected. The configuration panel on the left remains the same. The table on the right now displays the coefficient for the selected feature 'Néo'.

Feature	Coefficient	
80	Néo	0.3674

58

Dimensionality reduction: Intuition

Different goals, but the same principle – find the “good perspective”



PCA – increase variation
PLS-DA – increase class separation



“New perspective’s” coordinate axes
 = **Components**
 = linear combinations of features

<https://www.atmosera.com/blog/principal-component-analysis/>

59

59

sPLS-DA

Classification – prediction of categorical variable (binary or multinomial)

Assumptions:

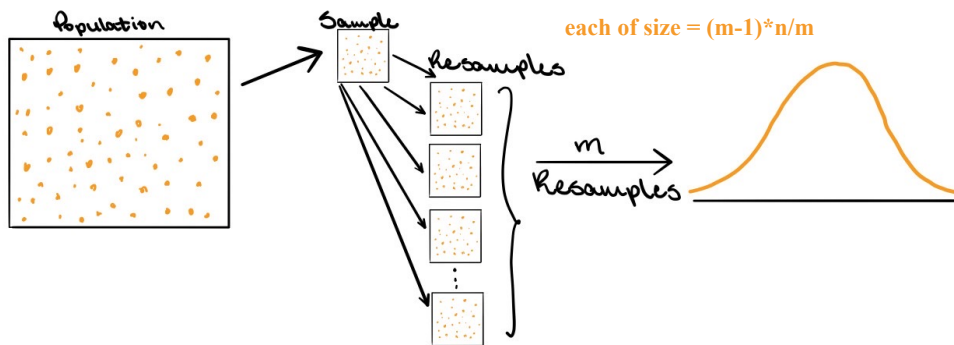
- **linearity** – classes can be separated by linear combinations of variables
- **no perfect multicollinearity**
- independence of observations
- no strong outliers
- no extremely imbalanced classes
- **only numeric predictors accepted**

60

60

Model Robustness

Brute force simulation approach: Bootstrapping or **Cross-Validation**



<https://www.linkedin.com/pulse/bootstrapping-statistics-what-why-its-used-trist-n-joseph>⁶¹

61

Feature Selector

Import Clean Filter Missing Values Odds Ratios GLM sPLS-DA Multi-Block sPLS-DA Final GLM Download Selection

Select the name of the discrete response variable:

DCD

Submit

Choose cut-offs for graphics

Feature stability threshold:

0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

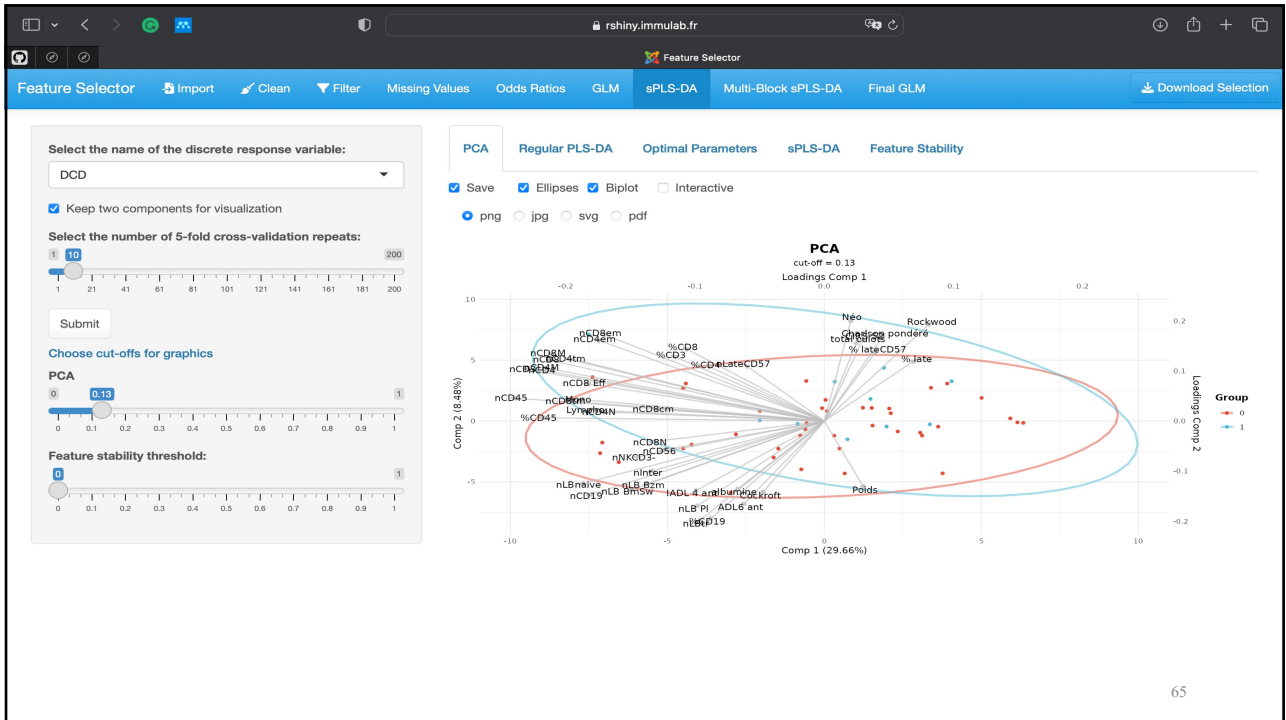
PCA Regular PLS-DA Optimal Parameters sPLS-DA Feature Stability

Save Ellipses Biplot Interactive

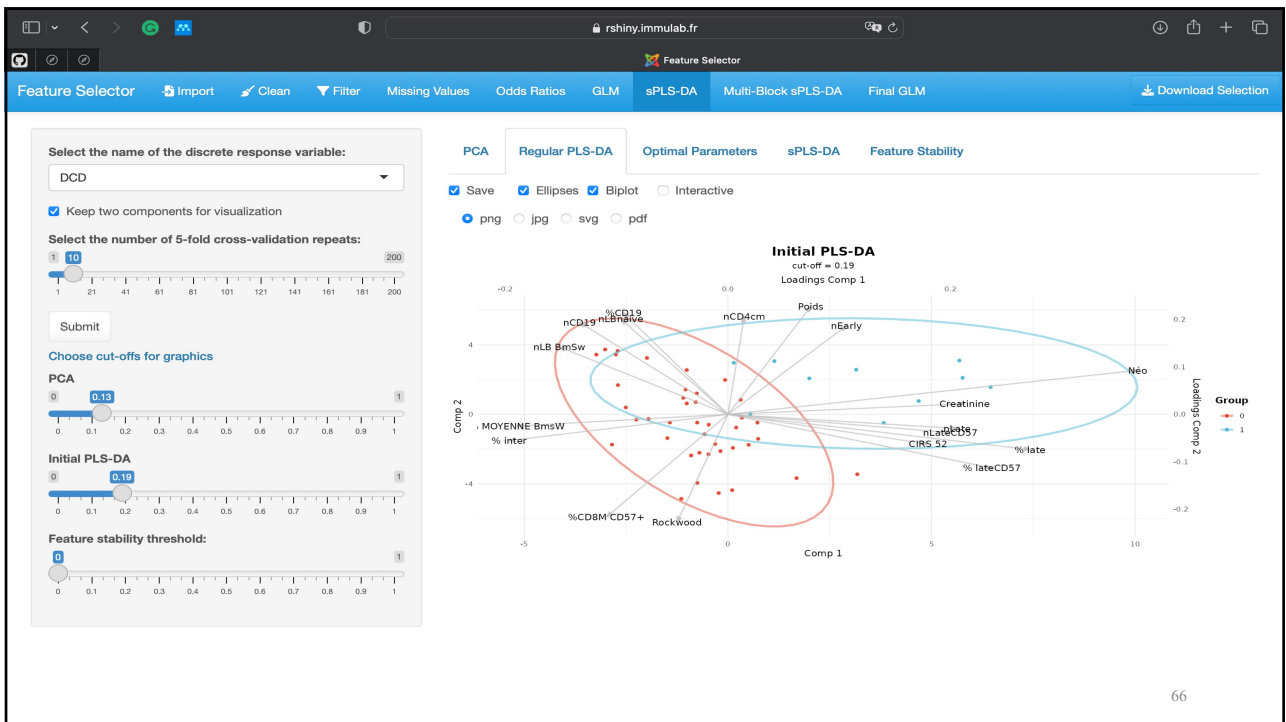
png jpg svg pdf

Be patient, this sPLS-DA takes some time...⁶²

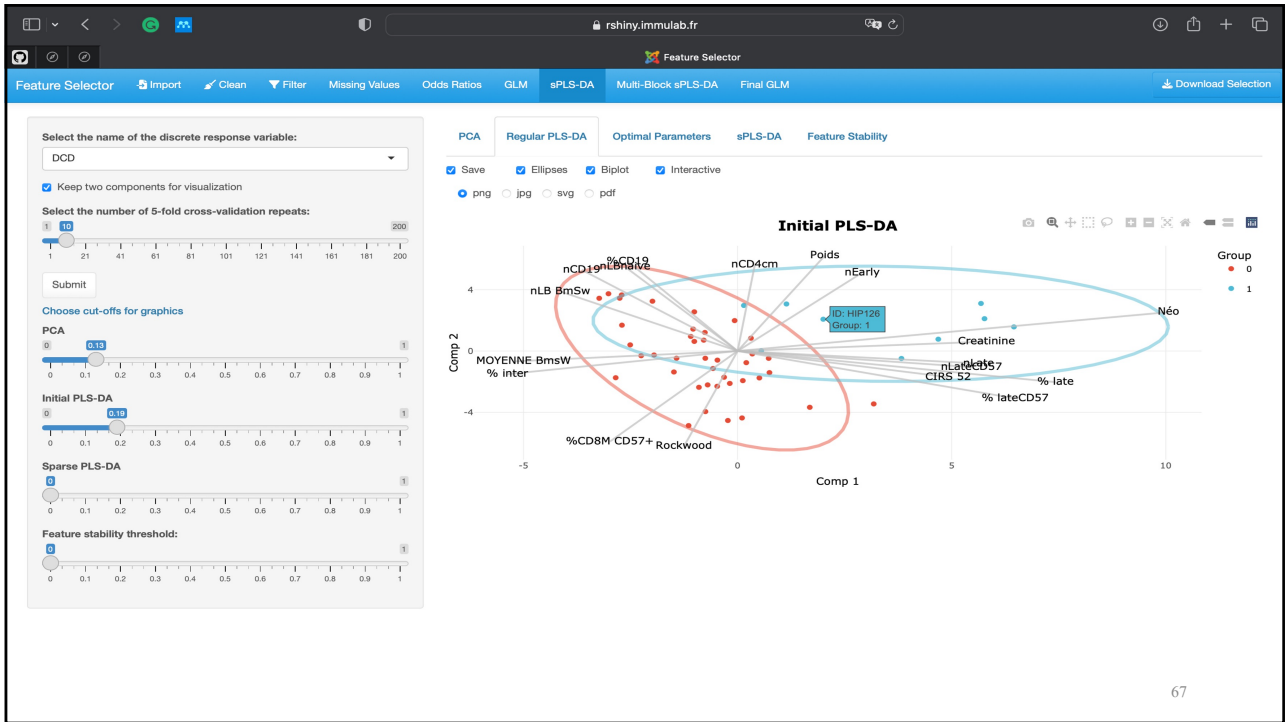
62



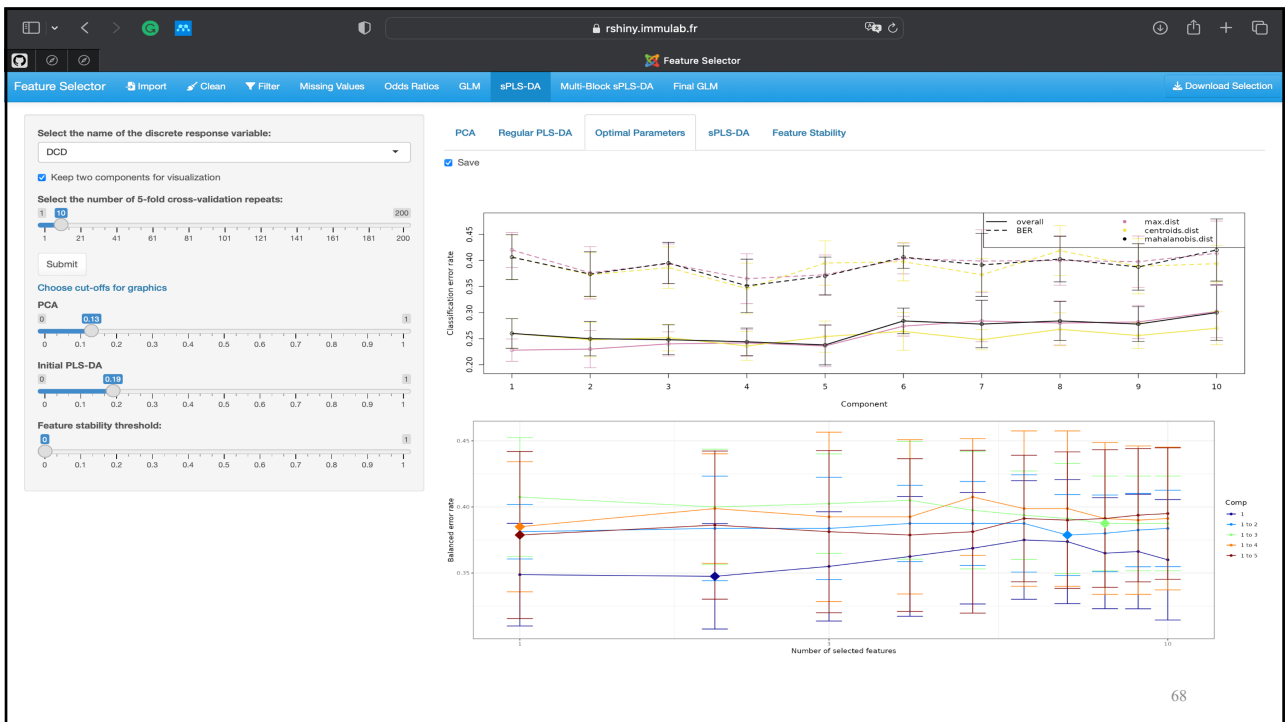
65



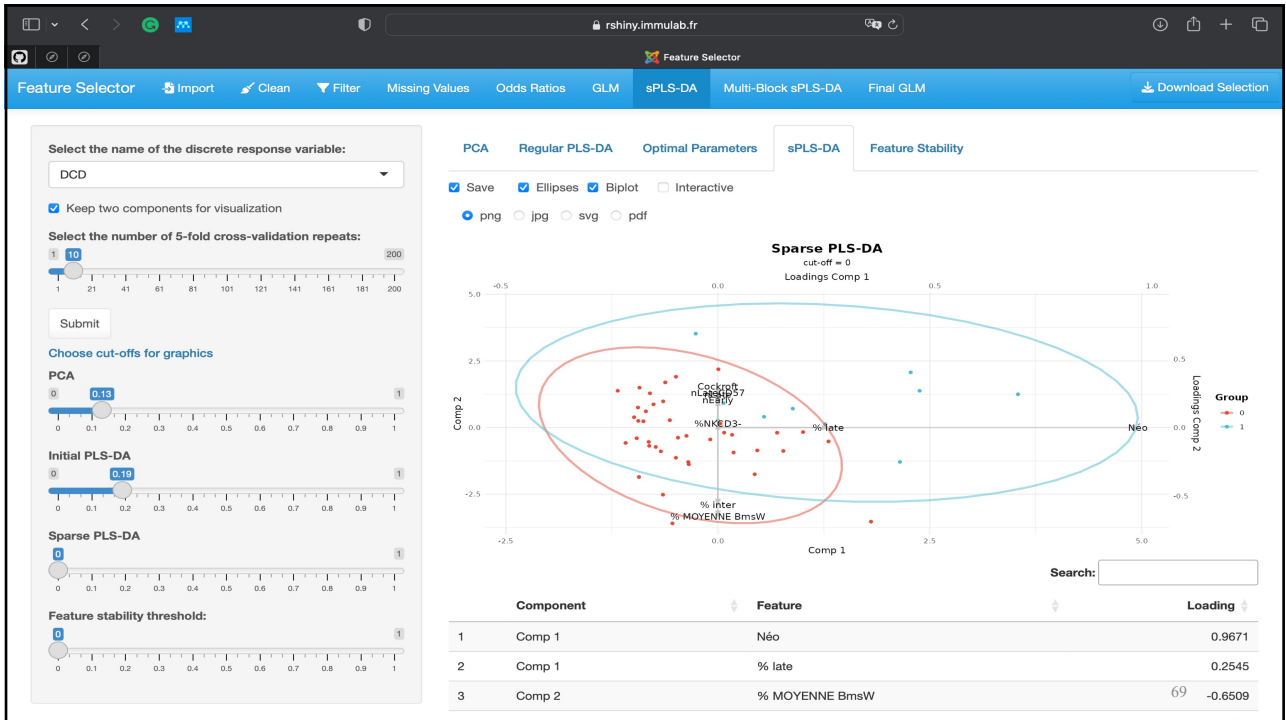
66



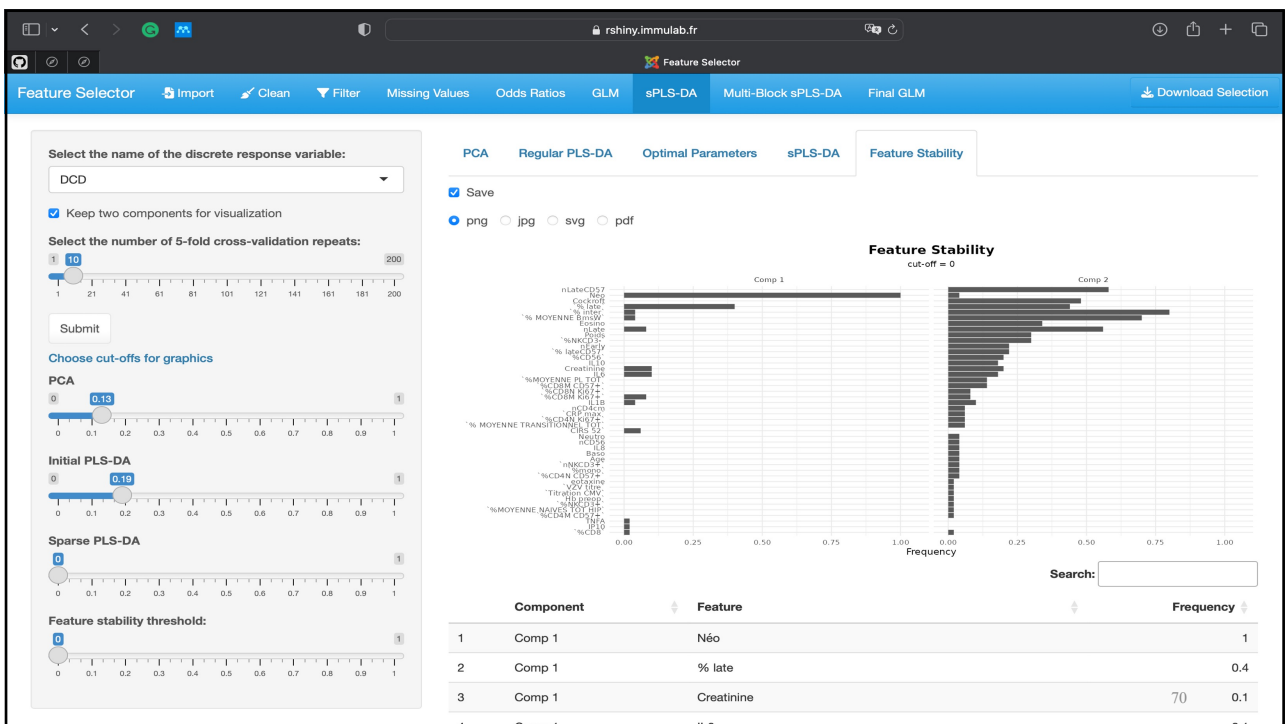
67



68



69



70

Multi-Block sPLS-DA (DIABLO)

- classification + integration of information from different blocks, takes into account inter-block interactions

Assumptions: in addition to sPLS-DA

- linearity– classes can be separated by linear combinations of variables inside and across blocks
- Contribution of variables from different blocks to the same component might mean the relation to the same phenomenon (e.g. pathological process, metabolic pathway)

71

71

Feature Selector | Import | Clean | Filter | Missing Values | Odds Ratios | GLM | sPLS-DA | Multi-Block sPLS-DA | Final GLM | Download Selection

Select the name of the discrete response variable:
DCD

Select the column indicating the block:
Categories

Submit

Choose cut-offs for graphics

PLS correlation circle:
0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

DIABLO:
0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

Feature stability threshold:
0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

Circos:
0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

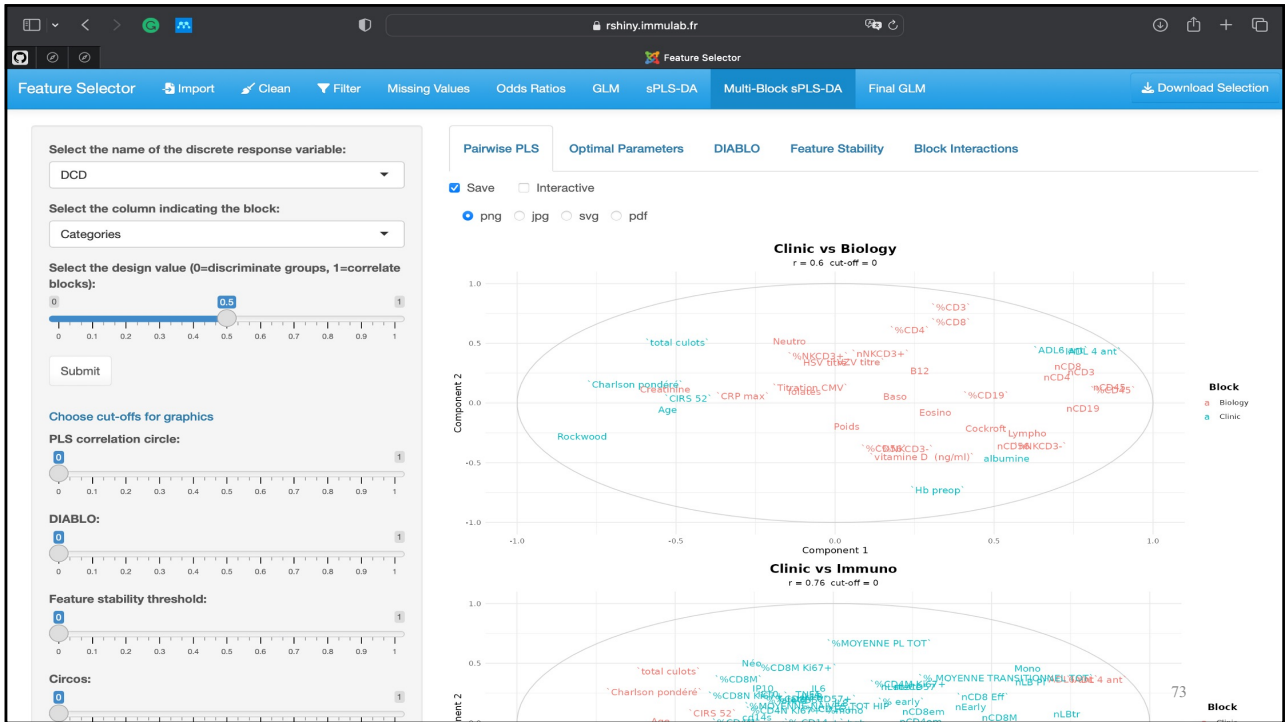
Pairwise PLS | Optimal Parameters | DIABLO | Feature Stability | Block Interactions

Save Interactive

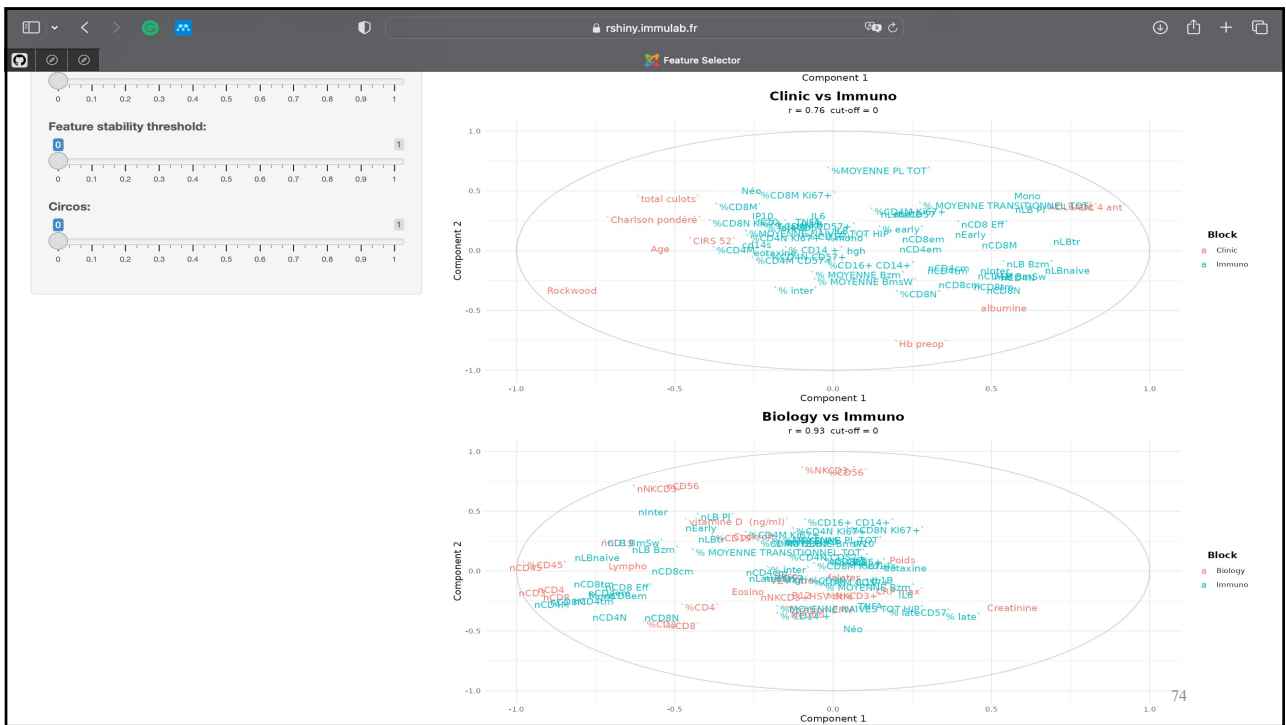
.png .jpg .svg .pdf

72

72



73

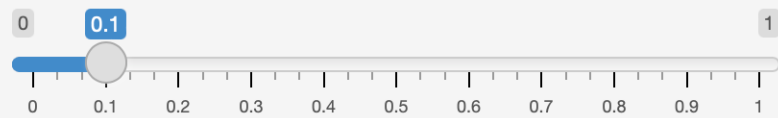


74

Choose the design value

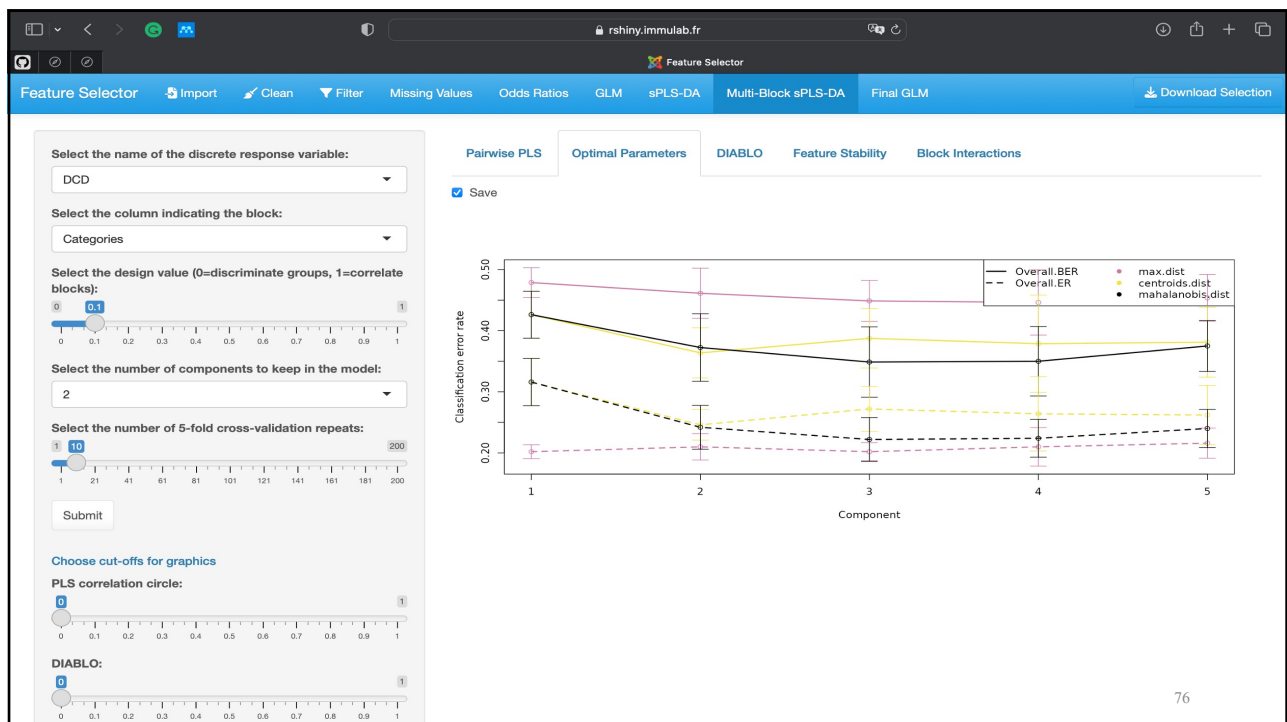
- Based on PLS results and the aim of the study
- From 0 to 1
- 0 - maximize discriminative ability
- 1 - maximize correlation between datasets

Select the design value (0=discriminate groups, 1=correlate blocks):



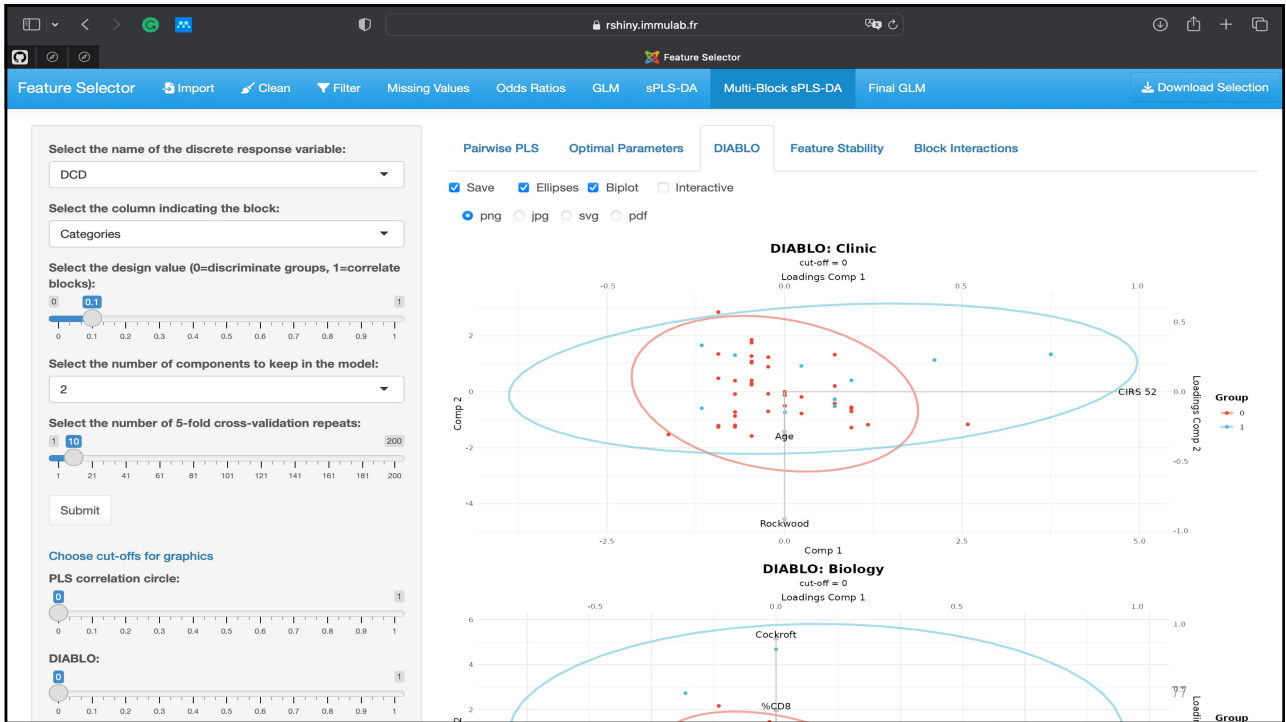
75

75

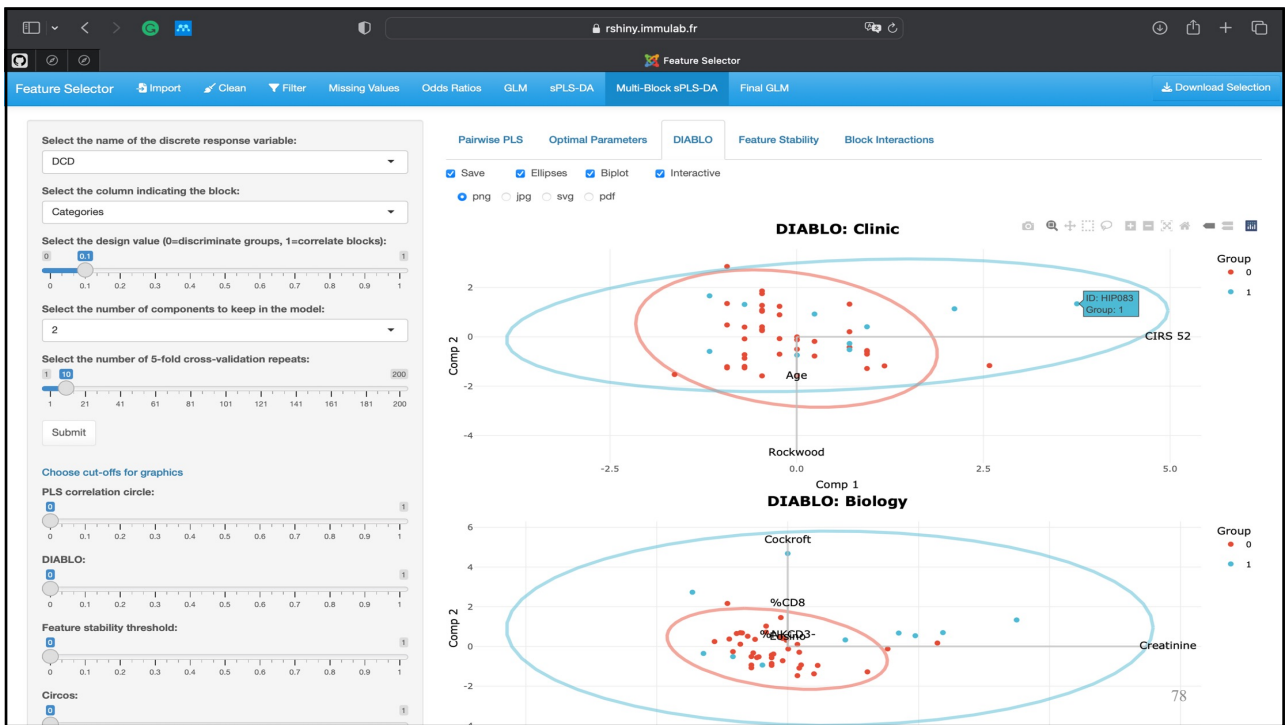


76

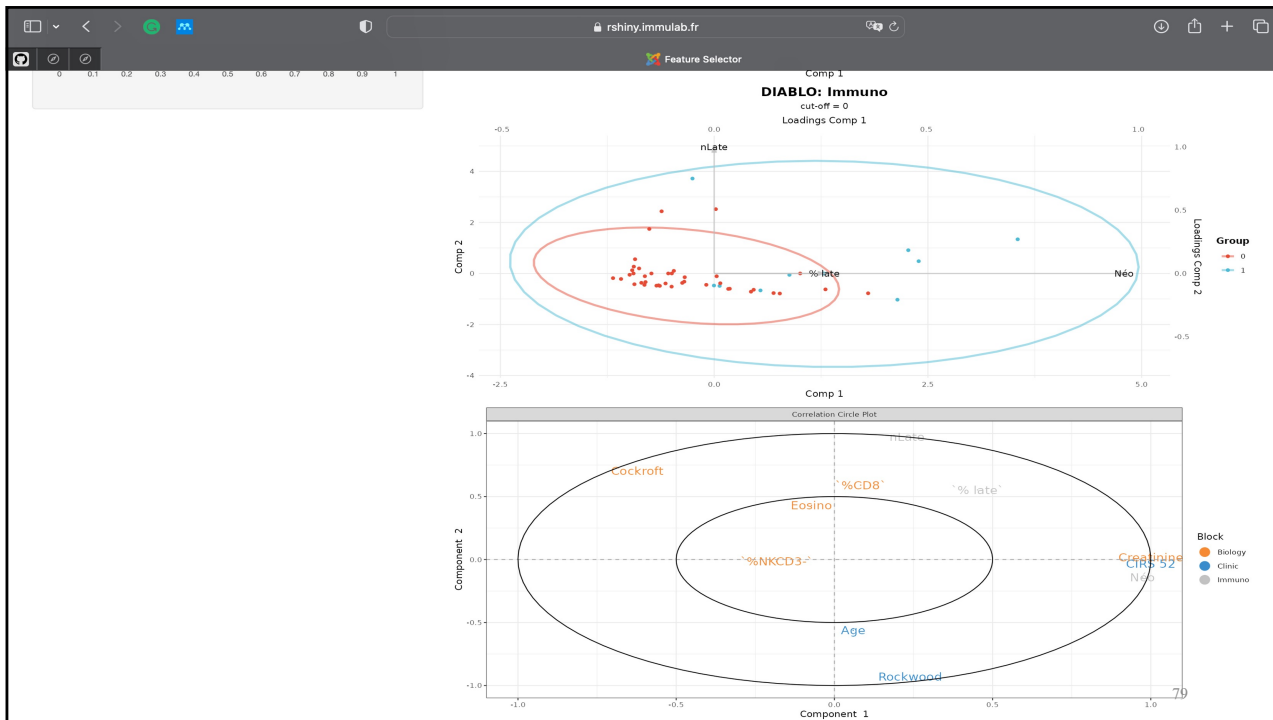
76



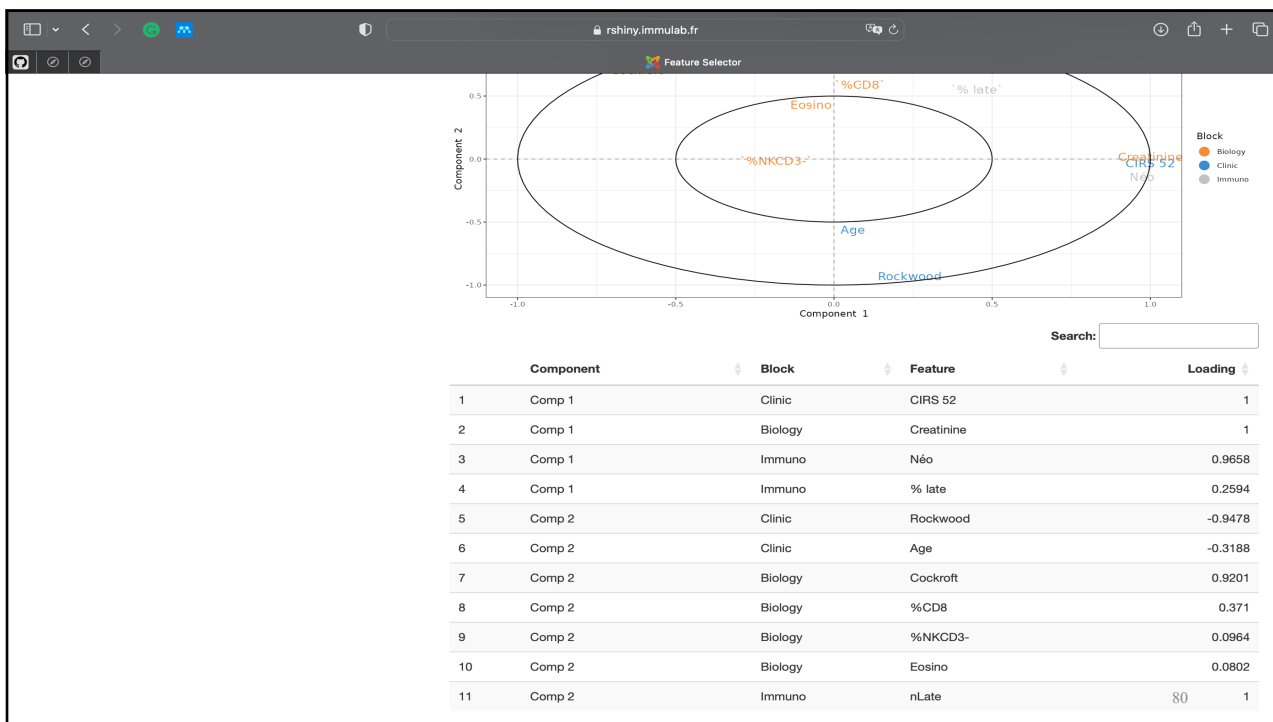
77



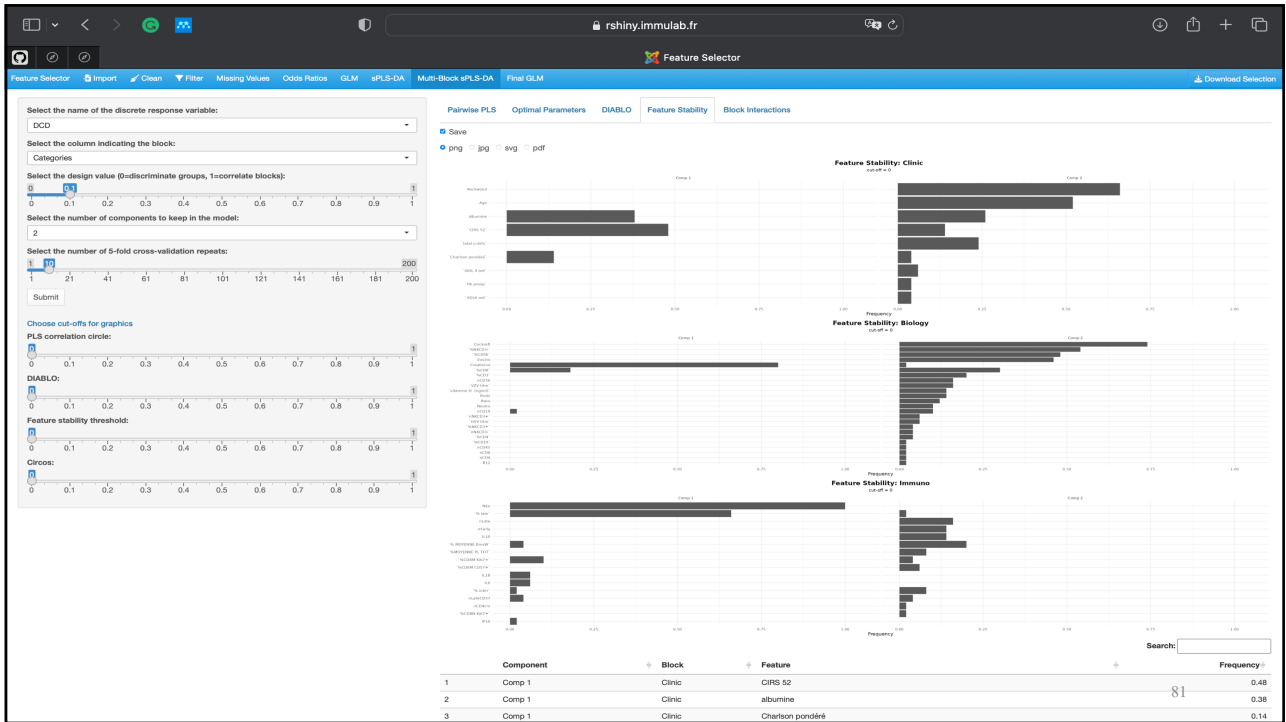
78



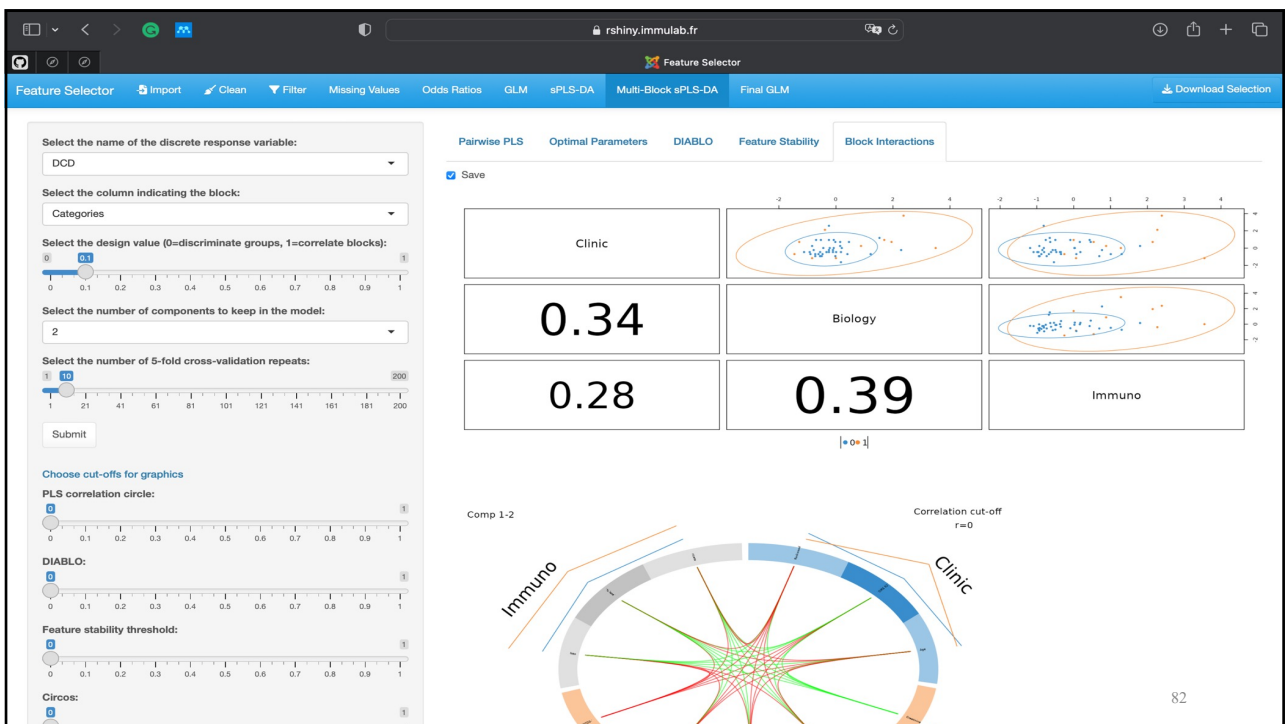
79



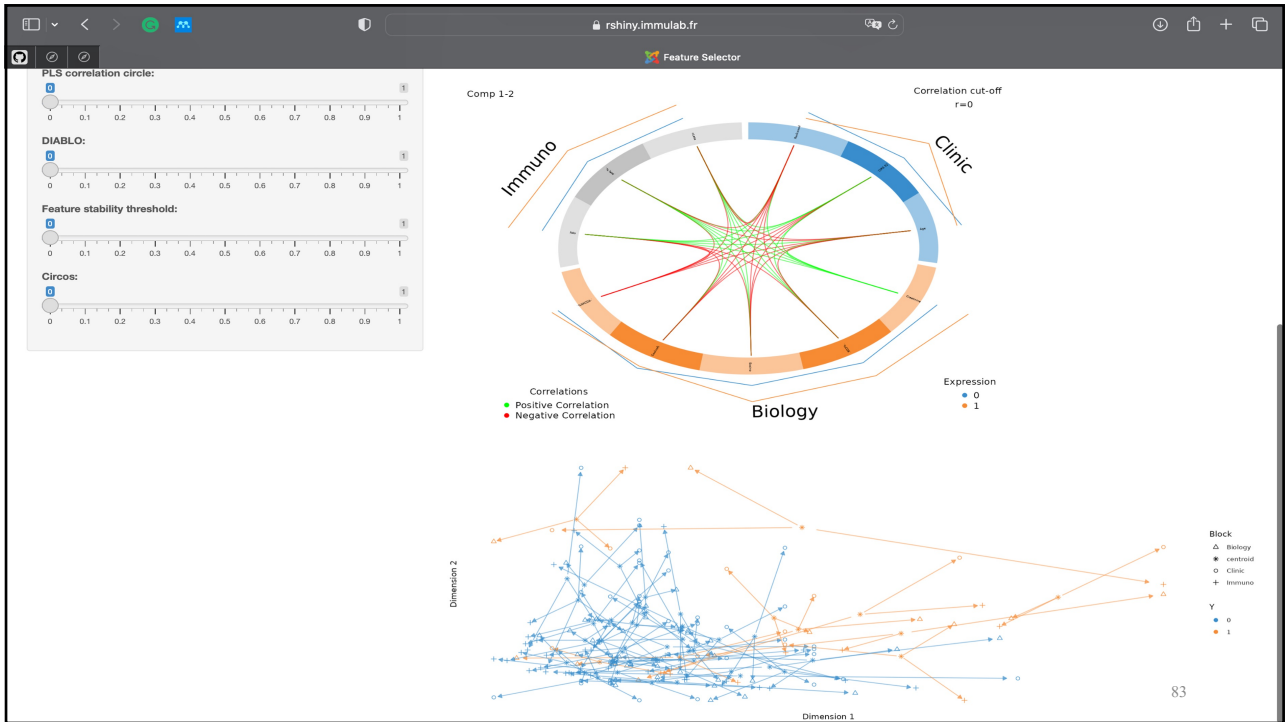
80



81



82

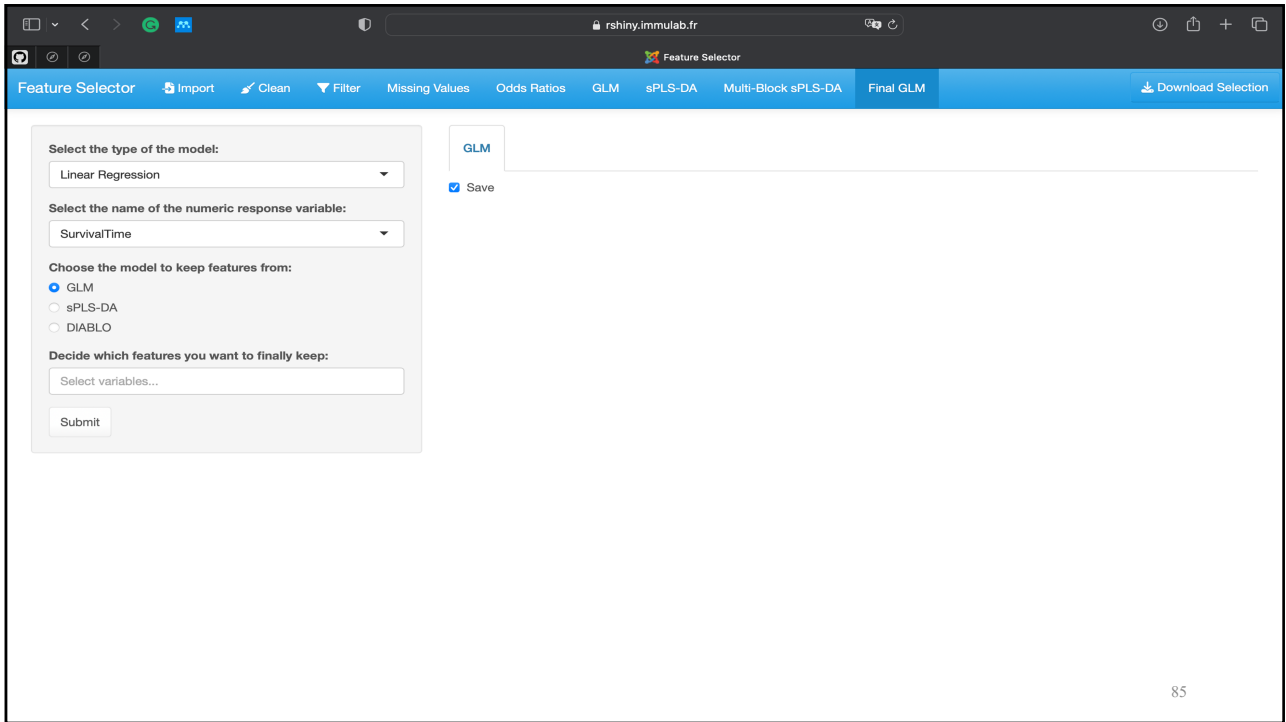


83

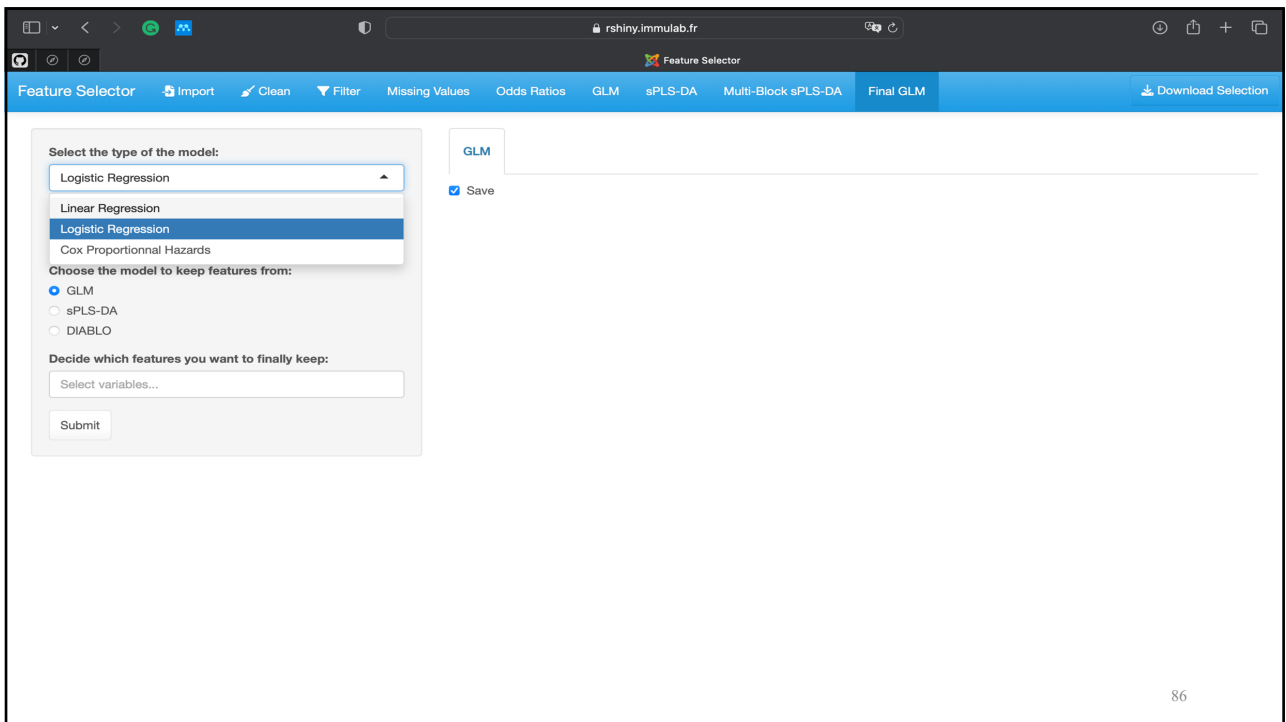
Final GLM

84

84



85



86

Feature Selector

Select the type of the model: Logistic Regression

Select the name of the binary response variable: DCD

Choose the model to keep features from:

- GLM
- sPLS-DA
- DIABLO

Decide which features you want to finally keep:

Select variables...

- Néo
- % MOYENNE BmsW
- Creatinine
- % late
- Poids
- %CD8M CD57+
- % inter
- Encl...

GLM

Save

87

87

Feature Selector

Select the type of the model: Logistic Regression

Select the name of the binary response variable: DCD

Choose the model to keep features from:

- GLM
- sPLS-DA
- DIABLO

Decide which features you want to finally keep:

Submit

GLM ROC

Save

DCD ~ Néo + `% MOYENNE BmsW` + `% late`

Call:
glm(formula = formula, family = model.type, data = data)

Deviance Residuals:

	Min	1Q	Median	3Q	Max
	-1.50620	-0.41443	-0.22937	-0.09025	1.91848

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-2.55751	1.49104	-1.715	0.08630 .
Néo	0.12876	0.04914	2.620	0.00879 **
% MOYENNE BmsW	-0.10082	0.06179	-1.632	0.10278 .
% late	0.06421	0.03655	1.757	0.07897 .

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 47.674 on 44 degrees of freedom
Residual deviance: 24.960 on 41 degrees of freedom
AIC: 32.96

Number of Fisher Scoring iterations: 6

Variance Inflation Factor (check for multicollinearity)

	Néo	% MOYENNE BmsW	% late
	1.494792	1.284764	1.219245

88

88

ROC curves

* in case of binary response

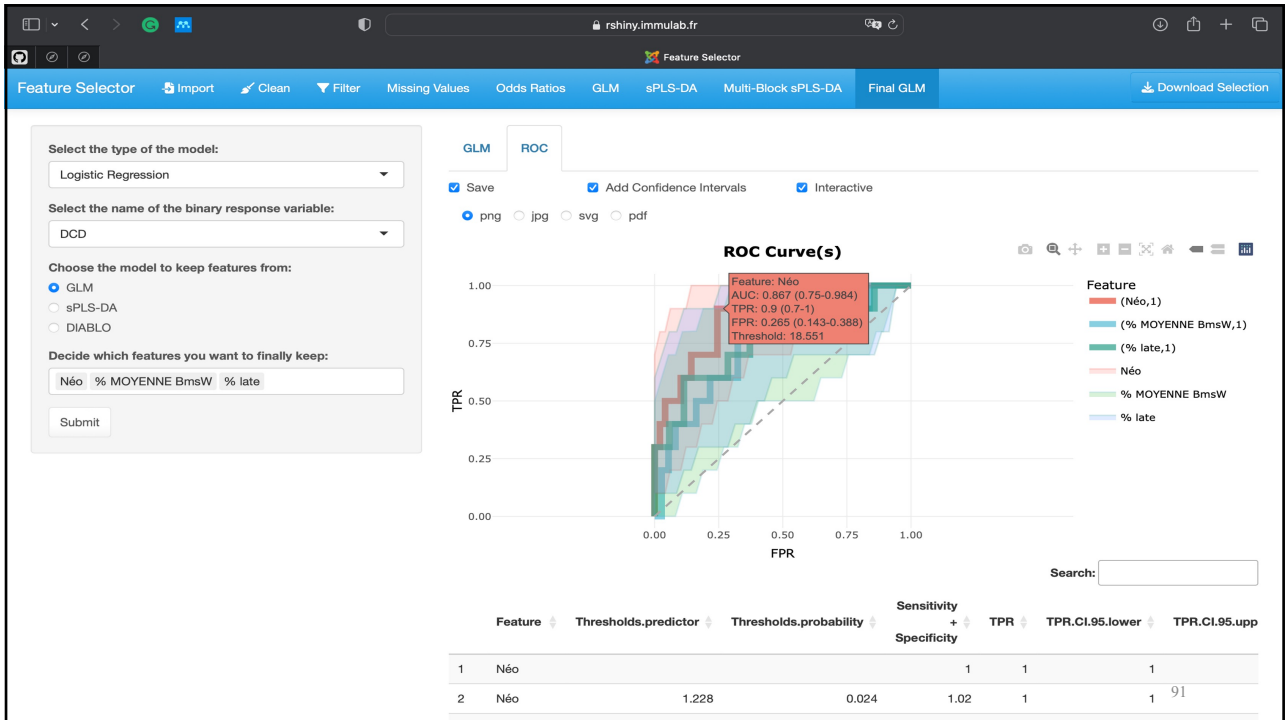
89

89

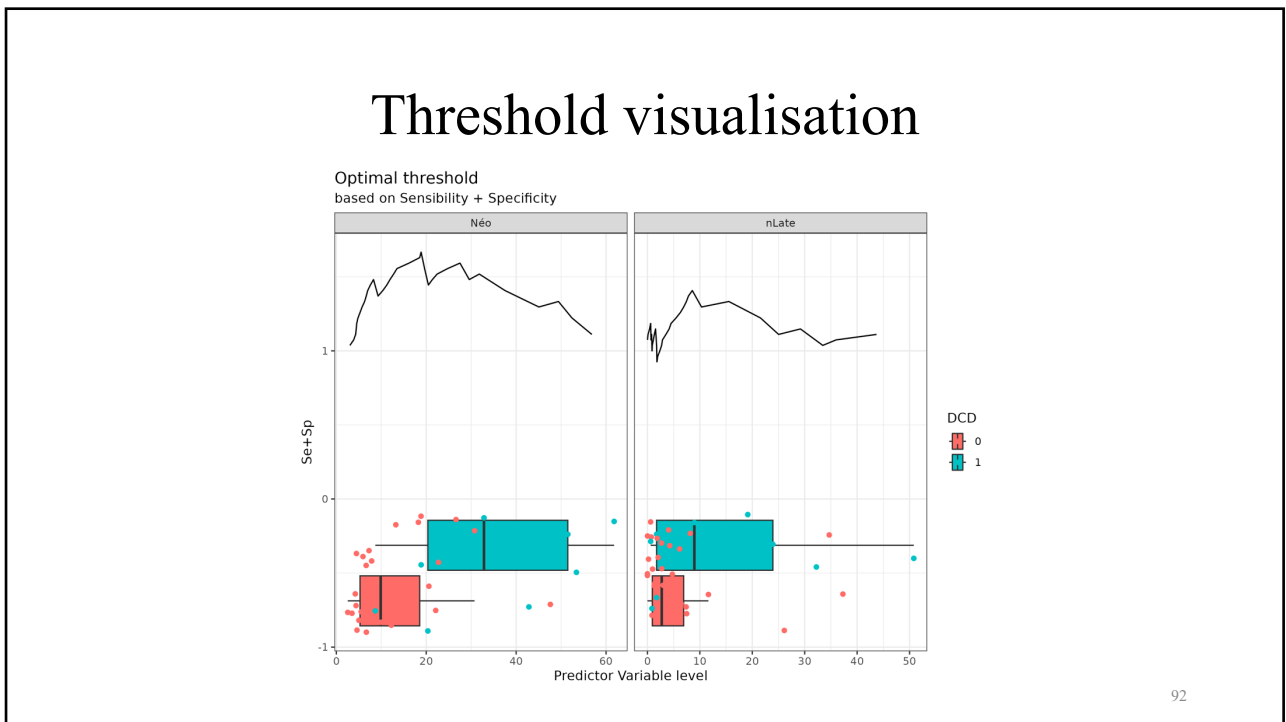
The screenshot shows the 'ROC' tab in the Feature Selector application. On the left, there are configuration options for the model (Logistic Regression), response variable (DCD), and feature selection criteria. The main area displays three ROC curves for features: 'Néo', '% MOYENNE BmsW', and '% late'. Below the plot is a table summarizing the performance of these features.

Feature	Thresholds.predictor	Thresholds.probability	Sensitivity + Specificity	TPR	TPR.CL95.lower	TPR.CL95.upper
1 Néo				1	1	1
2 Néo	1.228	0.024	1.02	1	1	90

90



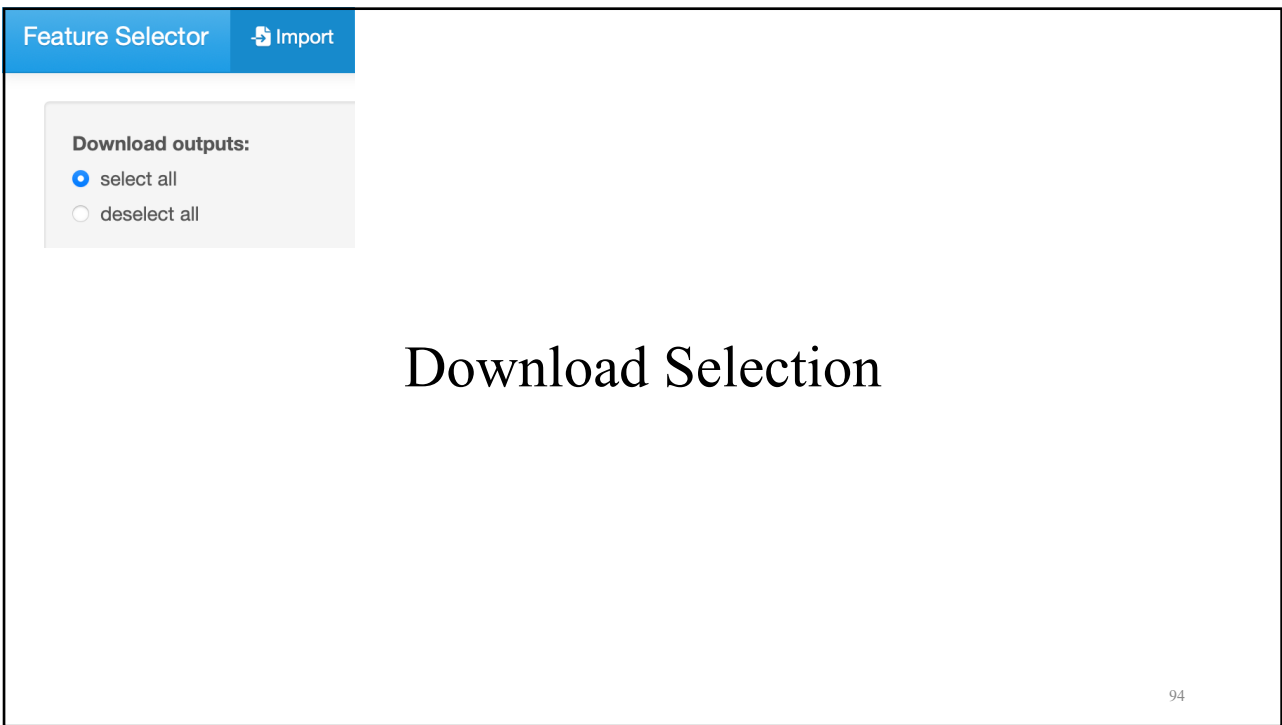
91



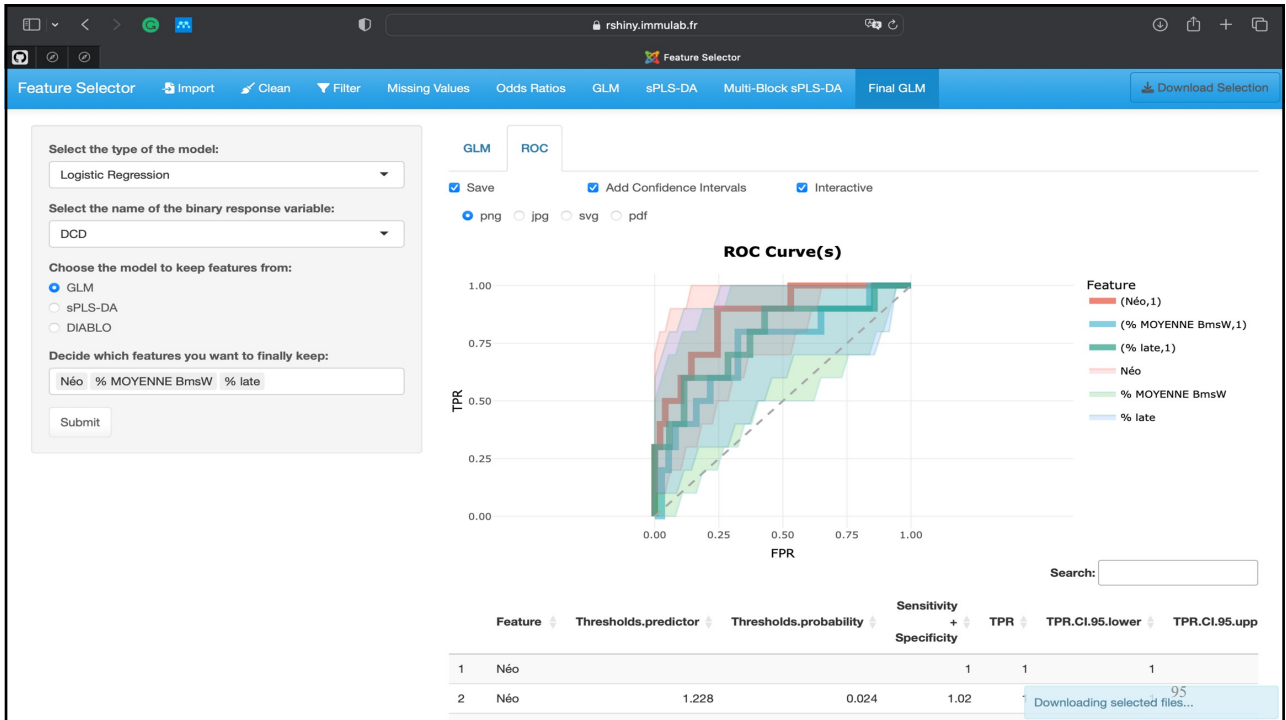
92



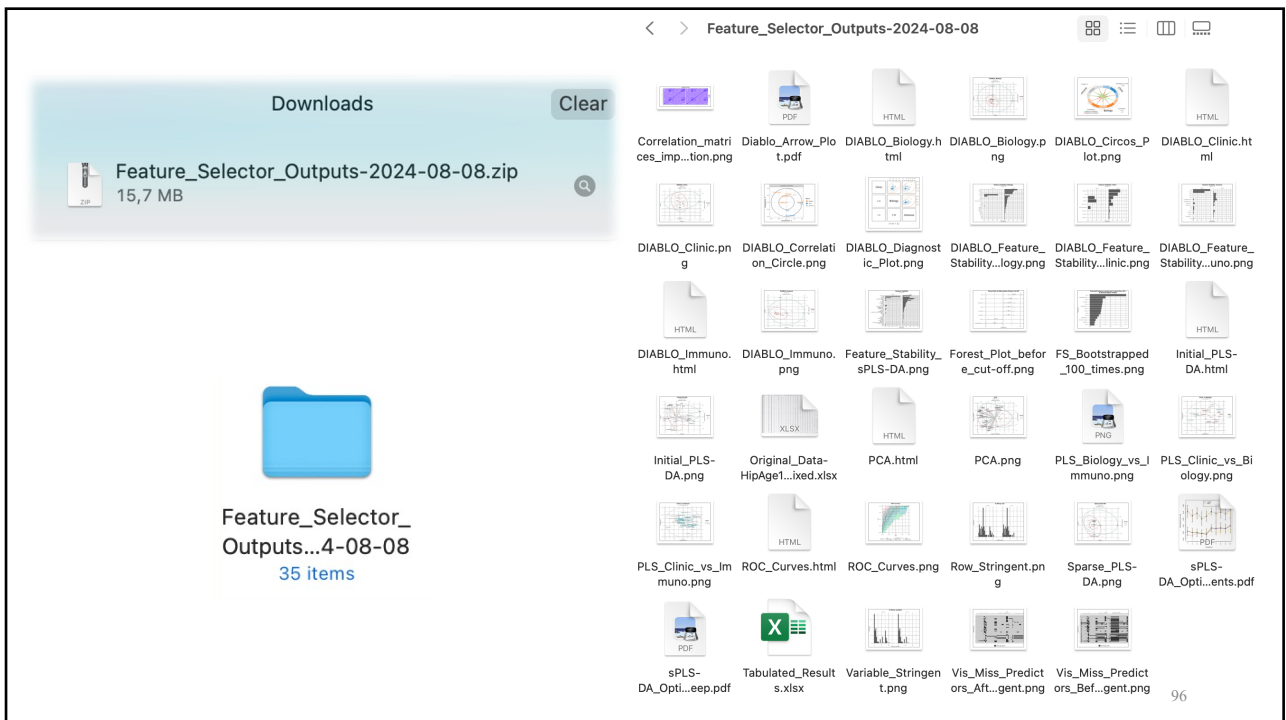
93



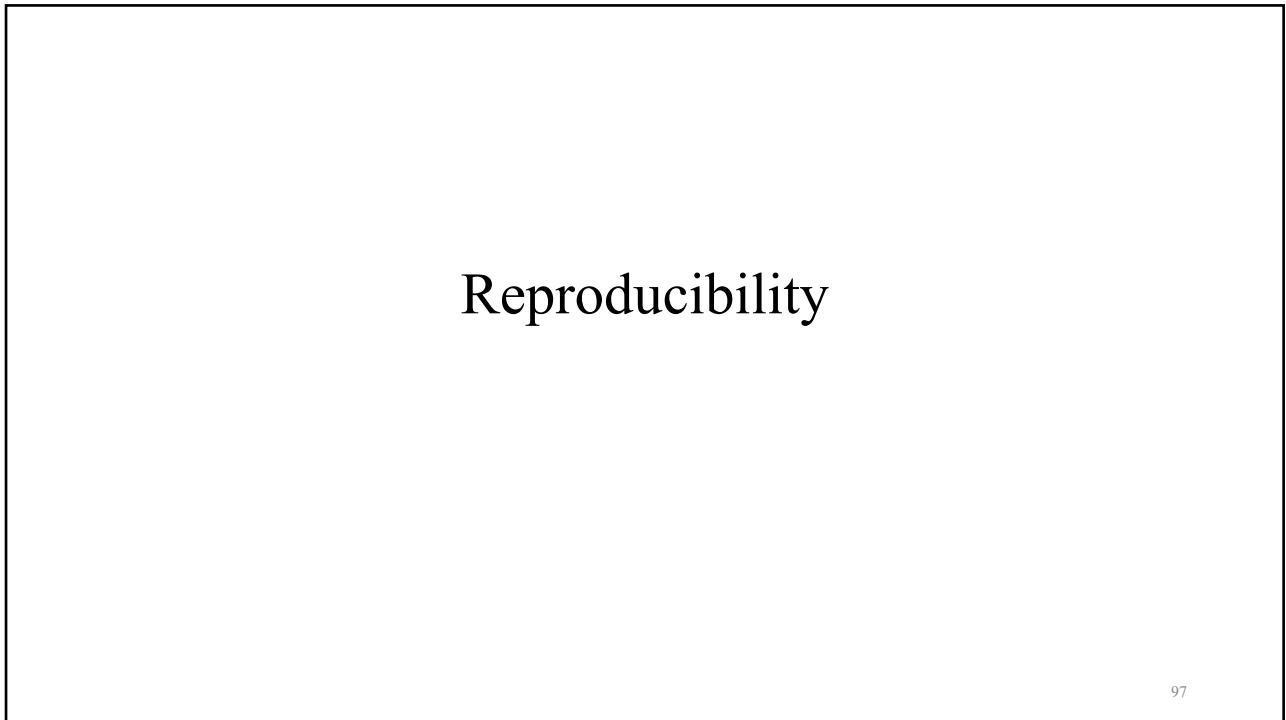
94



95



96



97

Parameter	Value	Comment
EXPORT	Source Data	
Source Data Sheet(s)		in case of several source sheets
Variable(s) to Merge by	Variable Information	
Variable Information Sheet	ColNames	
Variable Name Column		
CLEAN	5	for a variable that looks numeric, it will be considered as categorical if this value is below the threshold (threshold included)
Maximal Nb of Categories	5	for a variable that looks numeric, it will be considered as categorical if this value is below the threshold (threshold included)
Maximal Nb of Digits	5	for a variable that looks numeric, it will be considered as categorical if this value is below the threshold (threshold included)
FILTER		
ID Variable	PatientCode	
Response Variable(s)	DCD_SurvivalTime	excluded from predictors (I) if not kept explicitly with the next parameter
Response Variable(s) kept as Predictors		response variables used as predictors as well, e.g. prediction of NCD3 at D30 with NCD3 at D0
Response Filter Criteria		
Predictors Row Filter Criteria		
Predictors Column Filter Criteria		
Variables to Exclude		
Variables to Include Anyway		in spite of filters set above
MISSING VALUES		
Threshold as Proportion of NAs by Row	60.8	maximal tolerated % NA
Threshold as Proportion of NAs by Column	60.6	maximal tolerated % NA
Imputation Method	k Nearest Neighbors	
Central Tendency Method	5	if central tendency was selected
Nb of Neighbors	5	if kNN was selected
ODDS RATIOS		
Binary Response Variable	DCD	
GLM		
Model Type	Cox Proportional Hazards	
Response Variable	DCD	
"Survival Time" Variable	SurvivalTime	if Cox PH was selected
Imposed Alpha	5	0 = Ridge, 1 = Lasso
Optimal Alpha		chosen automatically if not imposed
Standardize Dummies	FALSE	
Nb of Bootstrapping Loops	100	
Oversampling	5	as the ratio of resampling size to the original sample size
sPLS-DA		
Factorial Response Variable	DCD	
Forced to Keep 2 Components	TRUE	for visualization, proposed if 2 is not optimal
Nb of Components	2	chosen by user among optimal values
Nb of Repeats of Cross-Validation	10	CV is 5-fold here, so multiply by 5 to get nb of models
Optimal Nb of Features for Each Component	2, 7	chosen automatically
DIABLO		= MultiBlock sPLS-DA
Factor Response Variable	DCD	
Variable Block Name	Categories	
Variable Blocks	Clinic, Biology, Immuno	
Design Value	0-1	maximize rather 0 = response discrimination, 1 = interblock correlation
Forced to Keep 2 Components	5	for visualization, proposed if 2 is not optimal
Nb of Components	2	chosen by user among optimal values
Nb of Repeats of Cross-Validation	10	CV is 5-fold here, so multiply by 5 to get nb of models
Optimal Nb of Features for Each Component by Block	c(1, 2), c(1, 4), c(1, 1)	chosen automatically
FINAL GLM		
Model Type	Logistic Regression	
Response Variable	DCD	
"Survival Time" Variable		if Cox PH was selected
Kept Feature(s)	N4s, % MOYENNE smw, % late	features that user has decided to keep in the end (not necessarily the best choice)

98

Perspectives

- Other imputation methods
- Basic visualization (e.g. boxplots)
- Interactive KM / Cox curves
- Set seed to reproduce randomness

- User-friendly tooltips and help windows

99

99

Let's try the app!

<https://www.immulab.fr/cms/index.php/team/tools/lab-tools/feature-selector>

100

100