

# CLASSIFICAZIONE del TUMORE DUTTALE PANCREATICO tramite BIOMARCATORI

UNIVERSITÀ degli STUDI di MILANO BICOCCA

Lombardi Mattia

**Abstract:** questo report intende analizzare, mediante classificazione supervisionata, il dataset aggiornato che era stato adoperato nello studio della Dott.ssa Debernardi per una previa valutazione del tumore pancreatico (Dec 2020, Journal PLOS Medicine) [[1]]. I metodi che verranno presentati sono il classificatore EDDA e un modello regressivo di misture finite gaussiane (GMM). Sono riportati, con classi bilanciate, individui: sani, aventi una malattia epatobiliare benigna o aventi un adenocarcinoma duttale pancreatico.

Il dataset è stato scaricato dalla piattaforma Kaggle:

<https://www.kaggle.com/datasets/johnjdavisiv/urinary-biomarkers-for-pancreatic-cancer?select=Debernardi+et+al+2020+data.csv>

## 1.INTRODUZIONE

L'adenocarcinoma duttale pancreatico (PDAC) è uno dei tumori più letali. Una volta diagnosticato, il tasso di sopravvivenza a 5 anni è inferiore del 10%. Tuttavia, se diagnosticato per tempo, le probabilità di sopravvivenza e la prognosi possono migliorare. Sfortunatamente molti casi di tumore al pancreas non mostrano sintomi fino a quando non raggiungono stadi più avanzati. Per questo motivo è importante individuare tecniche diagnostiche, non invasive, che possano valutare la presenza del tumore nei primi stadi. Infatti, questo dataset è costituito principalmente da biomarcatori delle urine e uno del plasma.

## 2.STRUTTURA DATI

Il dataset è composto da 590 osservazioni di 14 variabili. I dati sono stati raccolti da diverse fonti: Barts Pancreas Tissue Bank, University College London, University of Liverpool, Spanish National Cancer Research Center, Cambridge University Hospital, and University of Belgrade.

Le variabili *diagnosis* (31% control, 35% diagnosi benigna e 34% diagnosi tumore duttale pancreatico) e *sex* (49% M e 51% F) sono perfettamente bilanciate, questo suggerisce che

sia stato effettuato un *retrospective sampling*.

Quest'ultimo prevede di ottenere un dataset come concatenazione di sottocampioni, fin quando non si ottengono classi bilanciate da poter poi adoperare nel *training set*. Pertanto, il dataset va inteso per fini diagnostici e non descrittivi del fenomeno in questione.

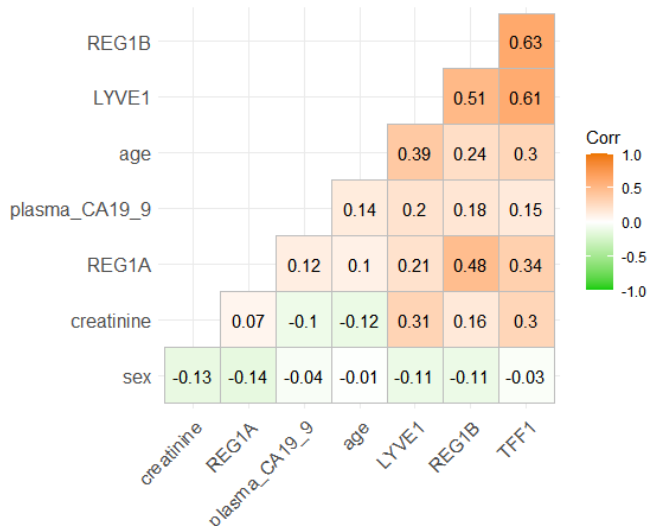
Una presentazione sintetica delle variabili coinvolte nel dataset:

variabili	descrizione
sample_id	stringa che identifica il paziente
patient_cohort	se l'osservazione proviene dal nuovo campione
sample_origin	da quale clinica viene il campione
age	età del paziente
sex	sexo del paziente (1-M;2-F)
diagnosis	(1-control,2-benigno,3-cancro)
stage	(IA, IB, IIA, IIIB, III, IV)
benign_sample	diagnosi per malattia benigna
plasma_CA19_9	anticorpo monoclonale che risulta spesso elevato nei pazienti con tumore pancreatico
creatinine	biomcatore della funzione dei reni
LYVE1	endotelio dei vasi linfatici
TFF1	"Trefol Factor 1", biomcatore coinvolto nella riparazione tessuti
REG1B	proteina associata alla rigenerazione del pancreas
REG1A	""", aggiunta successivamente per valutazione e confronto con REG1B

Purtroppo, le variabili *plasma\_CA19\_9* e *REG1A* presentano rispettivamente il 40.68% e 48.14% di dati mancanti, che rappresentano il 6.3% del totale. Dato che supera la canonica soglia del 5% non possono essere trascurati. Inoltre, le osservazioni mancanti risultano più frequenti nei pazienti senza diagnosi maligna, quindi i dati mancanti sono NMAR (probabilmente i pazienti non malati non eseguono di routine tali analisi). Pertanto, si è deciso di procedere sia tramite la non considerazione delle variabili (caso A), che tramite imputazione dei dati mancanti (caso B). In tal senso è stata adottata la PPCA (Probabilistic Principal Component Analysis), metodo che assume la normalità delle variabili e, perciò, ben si presta nell'imputazione dei dati mancanti in ambito medico-scientifico [[2]]. Inoltre, è un metodo robusto anche in caso di >15% *missings values*.

### 3. ANALISI ESPLORATIVA

L'analisi esplorativa è stata effettuata solo sui dati completi, in quanto si preferiva avere una visione più ampia possibile della relazione tra le variabili. Innanzitutto, visualizziamo la matrice di correlazione:

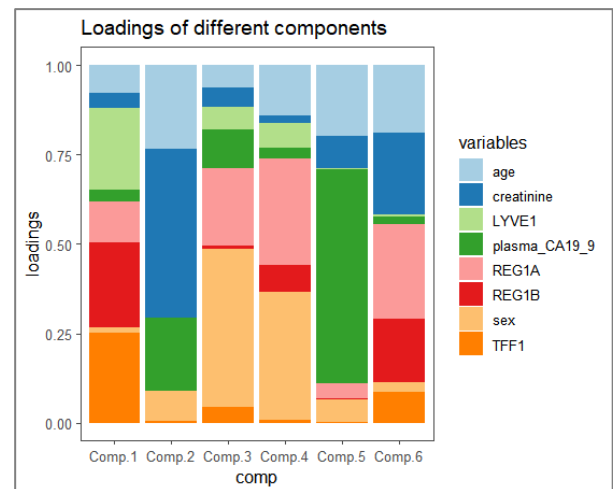


Come possiamo vedere non ci sono variabili particolarmente correlate tra loro ( $|p| > 0.80$ ). Pertanto, non vi è necessità di escludere alcuna variabile in quanto ridondante.

### PRINCIPAL COMPONENTS ANALYSIS (PCA)

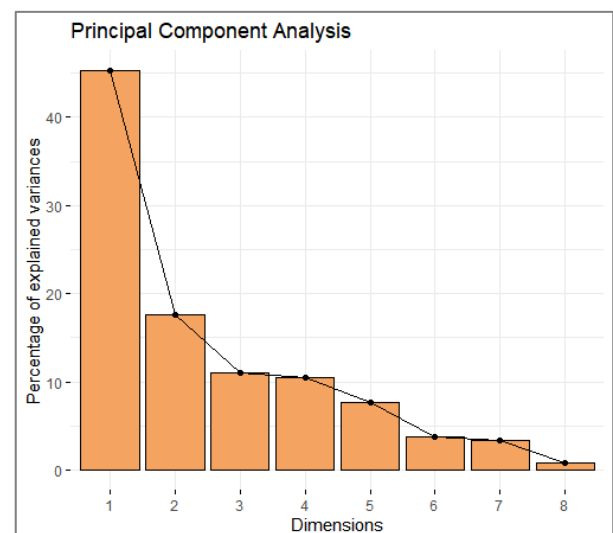
L'analisi delle componenti principali viene solitamente usata per ridurre la dimensionalità di un dataset. In questo caso, invece, si è preferito adoperarla più a titolo esplorativo. Infatti, tale scomposizione permette di comprendere le relazioni e la significatività delle componenti.

Le prime 6 componenti principali spiegano il 91% della variabilità interna. Pertanto, si è deciso di visualizzare i *loadings* di quest'ultime:

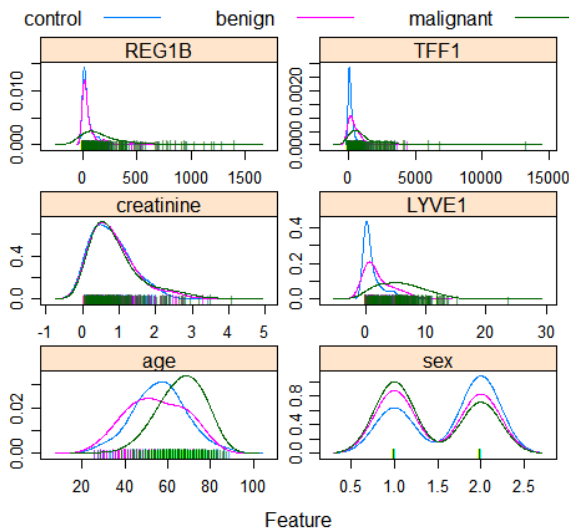


I loadings sono i coefficienti delle combinazioni lineari che definiscono le componenti principali. Quindi restituiscono una misura del contributo di ogni variabile alla componente considerata (1:6).

Si riporta, sottostante, anche il grafico che visualizza la percentuale di varianza spiegata per ogni componente:



Risulta importante anche valutare la distribuzione delle variabili. Come si può notare dal grafico sottostante i pazienti con diagnosi maligna hanno mediamente un'età più avanzata e presentano valori più elevati dei biomarcatori coinvolti nella rigenerazione del pancreas.



#### 4. EDDA CLASSIFIER

Il classificatore EDDA (Eigenvalue Decomposition Discriminant Analysis) è basato sulla scomposizione della matrice di varianza. Tale metodo di scomposizione viene detto VSO (Volume-Shape-Orientation) e permette la costruzione di una più flessibile struttura di variabilità:

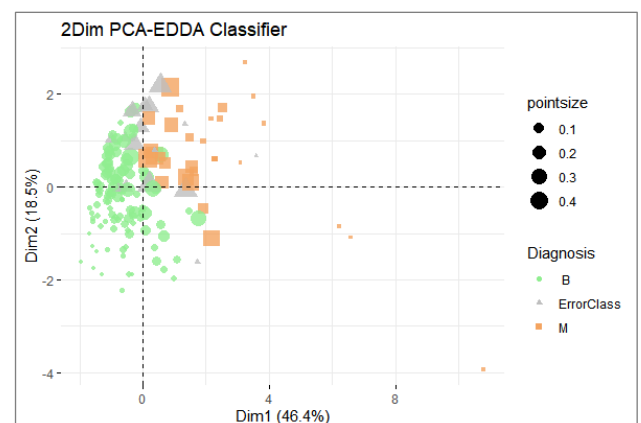
$$\Sigma_j = \lambda_j D_j A_j D_j'$$

Le analisi presentate riguarderanno il dataset che non presenta osservazioni imputate (caso A), in quanto l'aggiunta della covariate aggiuntive non aveva significativamente migliorato il modello. Inoltre, al fine di migliorare la *performance*, si è deciso di suddividere la popolazione in pazienti che non presentano metastasi (B) e pazienti aventi tumore maligno (M). Si procede suddividendo il dataset nella suddivisione training- test (75%-25%). Sono state effettuate 500 iterazioni per trovare il modello che minimizzasse il MER, tramite convalida incrociata (CV-10), e il BIC. Si sono poi ordinati i risultati mediante l'algoritmo *radix sort*. In tal senso il modello migliore è risultato *Gaussian\_pk\_L\_Ck* (**EVV**), ovvero un modello che presenta struttura

di variabilità con volume eguale, forma variabile ed orientamento variabile. Il secondo miglior modello è *Gaussian\_pk\_L\_C* (EEE), il quale coincide anche con la LDA (Linear Discriminant Analysis), un modello che presenta vincoli sulla matrice di varianza molto stringenti. Si è poi valutata la capacità previsiva del modello sul test-set, confrontando la vera diagnosi con quella stimata dal classificatore. Si è ottenuta la seguente matrice di confusione:

		Target		
		M	B	
Prediction	M	28	3	ACCURACY = 89.8%
	B	12	104	SENSITIVITY = 89.66%
				SPECIFICITY = 90.32%

Nonostante si abbia ottenuto una buona *accuracy* (89.8%), il classificatore non è riuscito a discriminare correttamente tutti i pazienti con metastasi. Purtroppo alcuni valori come la *TFF1* possono risultare molto elevati anche in caso di diagnosi secondarie (principalmente pancreatiti). Infatti si è ottenuto un ARI (Adjusted Rand Index) pari solo al 61%, il che indica che non era presente una netta differenza tra le due partizioni. Inoltre, l'informazione tra chi aveva una diagnosi secondaria benigna e chi era un soggetto sano è stata perduta. Nel grafico sottostante si rappresentano le osservazioni misclassificate considerando le prime due componenti principali. La grandezza dei punti rappresenta l'incertezza, infatti come possiamo notare le principali classificazioni errate erano contese tra le due nuvole di punti:



## 5. FINITE MIXTURES of REGRESSION MODEL (GMM)

Date le poche variabili a disposizione ed il fenomeno di studio, solo parzialmente noto, si è deciso di adottare una regressione di modelli mistura. Infatti, questi modelli riescono bene a catturare l'eterogeneità di una o più variabili non osservate (latenti). Legge distribuzione mistura:

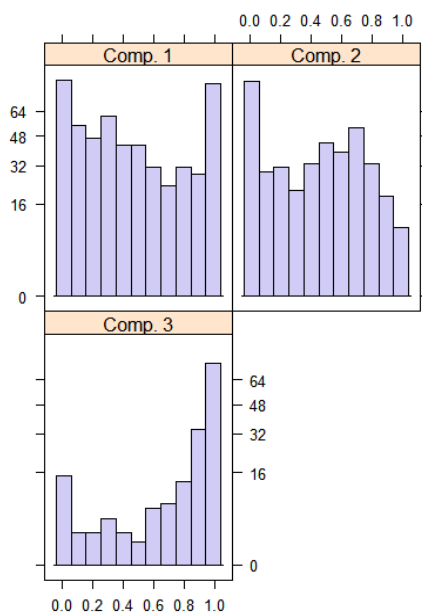
$$f(y | \varphi, x) = \sum_{k=1}^K \pi_k(x) \cdot f_k(y | \Theta_k, x)$$

In tal caso si considera la risposta  $y$  (diagnosi) distribuita normalmente e suddivisa nei  $k=3$  gruppi iniziali (C-B-M). Le covariate corrispondono alle variabili adoperate nella PCA con aggiunta delle osservazioni mancanti imputate (caso B), poiché in tal caso si è ottenuta una miglior classificazione.

*diagnosis~sex+age+creatinine+TFF1+LYVE1+plasma\_CA19\_9+REG1A+REG1B*

È stata utilizzata la libreria *flexmix* di Rstudio ed è stato implementato un Full MeM (Expert Network +Gating Network), ovvero sono stati modellizzati sia i parametri  $\varphi_k$  sia i pesi  $\pi_k$ . Inoltre il modello è stato scelto tramite 1000 iterazioni con il criterio di minimizzare sia EIC che BIC.

Rootogram of posterior probabilities > 1e-04



ICL	BIC	EIC
1184	925	0.62

Le iterazioni servono anche come meccanismo di controllo, poiché l'algoritmo E-M tende ad incorrere in massimi locali. Come possiamo vedere dal grafico delle probabilità a posteriori, le quali rappresentano il grado di incertezza con il quale sono state allocate ai gruppi le etichette, e dalla distanza Kullback-Leiber, le componenti B-M sono quelle tra loro più vicine e quindi più difficili da collocare.

Distanza K-L	
Comp(1-2 ; B-M)	2455.1
Comp(2-3 ; M-C)	150142
Comp(1-3 ; B-C)	75448.6

Gli *outliers* delle variabili *plasma\_CA19\_9* e *TFF1* sono stati rimossi in modo tale che il modello potesse differenziare maggiormente le componenti (rimangono 534 osservazioni). Qui di seguito visualizziamo la matrice di confusione ottenuta:

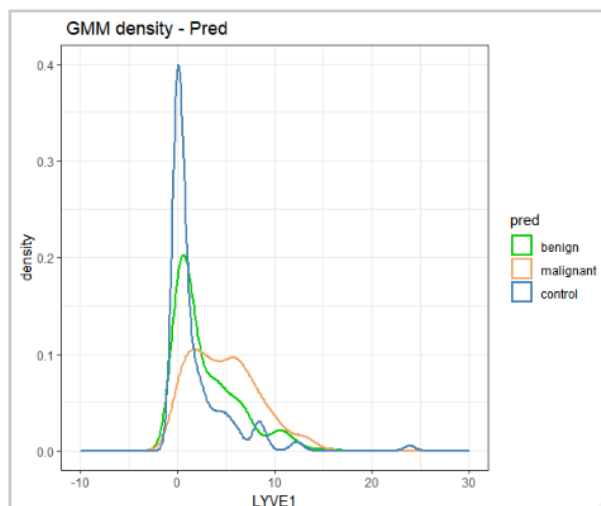
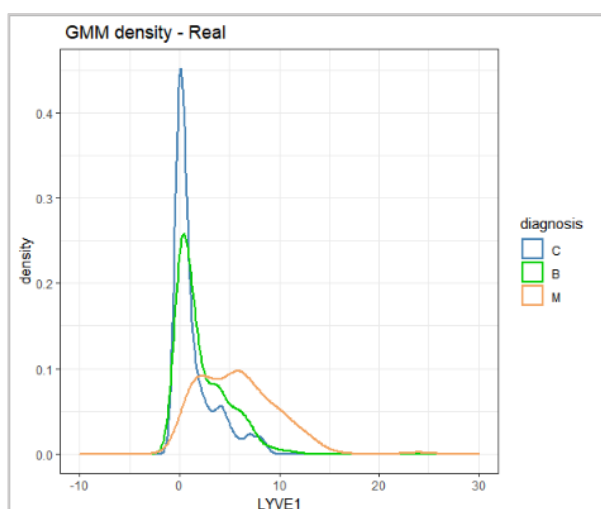
		Target			
		M	C	B	
Prediction	M	166		20	
	C	9	121	7	
	B	24	17	170	

ACCURACY = 85.58%  
SENSITIVITY(M) = 89.25%  
SPECIFICITY(M) = 90.52%

Gli errori di primo e secondo tipo sono stati riportati solamente per la diagnosi maligna (M), in quanto è la componente che principalmente ci interessa discriminare. Notare che nessun paziente del gruppo di controllo è stato classificato come malato, infatti come possiamo notare la distanza di KL è molto elevata tra M-C.

Dall'analisi della regressione otteniamo che: la componente C è stata modellizzata da *sex(\*\*\*)*, *age(\*\*\*)*, *LYVE1(\*\*\*)* e *REG1B(\*)*. Mentre la componente M: *intercept(\*\*\*)*, *age(\*\*)*, *TFF1(\*\*)*, *sex*, *plasma\_CA19\_9*, *LYVE1*, *REG1A*. La componente B, invece, solamente da: *intercept(\*\*\*)* e *LYVE1*. Gli asterischi corrispondono al livello di significatività dei coefficienti di regressione dei rispettivi modelli. (\*\*\*) - 0 ; \*\* - 0.01 ; \* - 0.05)

I modelli di misture gaussiane possono essere adoperati quando risulta difficile catalogare le osservazioni come appartenenti ad uno o all'altro gruppo, poiché riescono a classificare le unità in base al modello mistura, in tal caso normale, che meglio descrive ogni componente. Per valutare graficamente la bontà della regressione di misture gaussiane possiamo visualizzare il grafico della densità della variabile *LIVE1* per le etichette previste e poi confrontarla con la distribuzione di frequenza di quelle reali (lo scambio di etichette è dovuto all'algoritmo E-M ). Grafici sottostante:



## 6.CONCLUSIONI

I metodi *model-based* presentati sono risultati più che adeguati a prevedere e modellizzare, con una buona accuratezza e precisione, la diagnosi di PDAC (adenocarcinoma duttale pancreatico).

Inoltre, in tal caso non avevamo a disposizione un grande numero di osservazioni (590), quindi tecniche non parametriche avrebbero potuto dare performance inferiori, oltre a non poterci restituire un'interpretazione della significatività delle variabili.

In realtà anche il classificatore EDDA avrebbe giovato di un dataset più grande e completo, in quanto avrebbe avuto un più ampio training set. Pertanto risulterebbe utile continuare a collezionare ulteriori dati riguardanti la patologia, perché ripetere le analisi su un dataset più fornito potrebbe migliorare di molto le performance.

Purtroppo raggiungere un'*accuracy* pari al 100% risulterebbe difficile in ogni caso, poiché alcuni malati di PDAC non presentano analisi dei biomarcatori notevolmente alterate rispetto ai pazienti sani o con diagnosi secondaria.

In tal senso approfondire con altre libraries, o metodologie di convergenza per l'algoritmo E-M, le regressioni di misture finite gaussiane (GMM) potrebbe risultare interessante, in quanto sono un modello che riesce a cogliere la presenza di variabili latenti. Pertanto potrebbero risultare particolarmente utili in questo specifico caso. Inoltre adoperare modelli statistici più avanzati, come quelli Bayesiani, potrebbe risultare interessante e dare notevoli spunti ai fini della classificazione.

## REFERENCES

[[1]]

<https://journals.plos.org/plosmedicine/article?id=10.1371/journal.pmed.1003489>

[[2]]

<https://www.sciencedirect.com/science/article/pii/S2352914819302783>

*In dedica e memoria di mio Nonno Luciano*

# CODICE ADOPERATO per GRAFICI ed ANALISI

(software adoperato: Rstudio)

mat: 837193

```

1 #Mattia Lombardi#
2
3 #project of Computational Statistics:
4 #Model based Classification of pancreatic cancer
5
6 #libraries:
7 library(tidyverse)
8 library(ggplot2)
9 library(mclust)
10 library(visdat) #missed data visualization
11 library(ggcorrplot)
12 library(GGally)
13 library(ggthemes)
14 library(factoextra) #pca library
15
16 #-----
17 #PCA and Exploration Data Analysis (EDA)
18 #-----
19
20 setwd("C:/Users/39349/OneDrive/Desktop/UNIMIB/dataset")
21
22 data<-read.csv("pancreatic_cancer.csv")
23 #visualize data
24 str(data)
25 view(data)
26
27 #cambio in numerica la variabile sex M-F
28 data$sex<-ifelse(data$sex=="M",1,2)
29 table(data$diagnosis)/590
30 table(data$sex)/590
31 #proportion are ok
32
33 #missing data
34 anyNA(data) #TRUE
35 vis_miss(data)
36
37 #purtroppo le variabili REG1A e plasma_CA19_9
38 #hanno molti valori mancanti
39 #per fare un analisi preliminare delle PCA
40 #consideriamo solamente le 209 osservazioni complete
41
42 data_comp<-data[complete.cases(data),]
43 #considero variabili solo numeriche
44 data_num<-data_comp[, -c(1,2,3,6,7,8)]
45 data_lab<-as.factor(data_comp[,6])
46 table(data_lab)/209
47
48 #possiamo notare che i missing values provenivano
49 #prevalentemente dal gruppo di controllo o con diagnosi benigna
50 #dati mancanti NMAR
51 str(data_num)
52
53 #correlazione
54 ggcorrplot(cor(data_num), hc.order = TRUE, type = "lower", lab = TRUE,
55             colors = c("green3", "white", "darkorange2"))
56

```



```

57 #PCA
58 pca<-princomp(data_num,cor=T)
59 pca
60
61 #rapp. della % di varianza spiegata dalle components
62 fviz_eig(pca, barfill = "sandybrown",barcolor="black",main = "Principal Component Analysis")
63 sum(pca$sdev[1:6]^2)/8
64 #le prime 6 components spiegano 91% varianza within
65
66 #grafici non presenti nel report:
67 #1)rapp. del grafico rappresentante le 2 componenti principali suddivise per sesso
68 groups <- as.factor(ifelse(data_num$sex==1,"M","F"))
69
70 fviz_pca_ind(pca,
71             col.ind = groups,
72             palette = c("#FC4E07","#00AFBB" ),
73             addEllipses = TRUE,
74             legend.title = "Sex",
75             repel = TRUE,
76             title=c("2Dim PCA-Sex")
77 )
78
79 #2)rappresentazione grafica 2PCA dimensions - creatinine
80 fviz_pca_ind(pca, col.ind=data_num$creatinine,
81             legend.title = "Creatinine",
82             repel = TRUE,
83             title=c("2Dim PCA-Creatinine")
84 )+scale_color_gradient2(low="white", mid="slateblue",
85                       high="orange", midpoint=0.6)
86
87 #costruisco il grafico dei loadings
88 scores<-pca$loadings[,1:6]^2
89 scores
90 data_scores<-as.vector(scores)
91
92 comp<- c ( rep("Comp.1",8),rep("Comp.2",8),rep("Comp.3",8),rep("Comp.4",8),rep("Comp.5",8),rep("Comp.6",8))
93 variables<-c("age","sex","plasma_CA19_9","creatinine","LYVE1","REG1B","TFF1","REG1A")
94 data_loadings<-data.frame(comp,variables,data_scores)
95 data_loadings
96
97 #grafico loadings
98 ggplot(data_loadings, aes(fill=variables, y=data_scores, x=comp)) +
99   geom_bar(position="stack", stat="identity")+
100   labs(title = "Loadings of different components",y="loadings")+
101   scale_fill_brewer(palette = "Paired")+
102   theme_few()
103 data_num
104 #mostro le funz di densità delle variabili
105 #suddivise per ogni componente
106 library(caret)
107 levels(data_lab)<-c("control","benign","malignant")
108
109 featurePlot(data_num[, -c(3,8)],data_lab,plot="density",
110            scales=list(x = list(relation="free"),y = list(relation="free")),
111            adjust = 1.5,pch="|",auto.key = list(columns = 3))

```

```

113 #-----
114 #EDDA CLASSIFIER
115 #-----
116
117 #rileggo i dati per evitare eventuali errori
118 data<-read.csv("pancreatic_cancer.csv")
119
120
121 #data manipulation:
122 data$sex<-ifelse(data$sex=="M",1,2)
123
124 #modifico le etichette in diagnosi senza metastasi e con cancro
125 data$diagnosis<-ifelse(data$diagnosis %in% c(1,2),1,2)
126
127 #considero solo le variabili complete e numeriche
128 dada<-data[, -c(1,2,3,6,7,8,9,14)]
129 #etichette reali
130 data_lab<-as.factor(data$diagnosis)
131
132 #per ottenere sempre stessi risultati
133 set.seed(123)
134
135 #training e test set
136 test_lab = sample(1: 590 ,trunc(0.25* 590),replace = F)
137
138 train<-as.data.frame(dada[-test_lab,])
139 data_train_lab<-data_lab[-test_lab]
140
141 #dataframe per inserire risultati
142 results = as.data.frame ( matrix ( nrow = 300 , ncol = 2))
143 colnames(results)=c("model","cv")
144
145 library(Rmixmod)
146 clas.train=mixmodLearn(train,data_train_lab,
147                        models = mixmodGaussianModel(family="all",equal.proportions=FALSE),
148                        criterion="BIC")
149 clas.train
150
151 for ( i in 1:300){
152   clas.train=mixmodLearn(train,data_train_lab,
153                          models = mixmodGaussianModel(family="all",equal.proportions=FALSE),
154                          criterion=c("CV","BIC"))
155
156   model=clas.train@bestResult@model
157   cv=clas.train@bestResult@criterionValue[1]
158   results[i,]=c(model,cv)
159 }
160
161 #ordino i risultati
162 results$cv = as.numeric ( results$cv )
163 results<-results[order(results$cv,method="radix"),]
164
165 #seleziono miglior modello
166 best.mod = mixmodLearn (train,data_train_lab,
167                        models = mixmodGaussianModel(listModels=results$model[1]))
168 best.mod["bestResult"]

```



```

169
170 #prediction
171 pred.values=mixmodPredict(data=dada[test_lab,],classificationRule=best.mod["bestResult"])
172 |
173 uncertainty=apply(pred.values@proba,1,min)
174 types=factor(pred.values@partition,levels =c(1,2),labels = c("B","M"))
175
176 act = factor ( data_lab [ test_lab ] , levels = c ("1","2") ,
177               labels = c (" B ", "M"))
178
179 prd = factor ( pred.values@partition , levels = c ("1","2") ,
180               labels = c (" B ", "M"))
181
182
183 library(caret)
184 conf.matrix = confusionMatrix ( data = act , reference = prd )
185 conf.matrix
186
187 tab_conf<-as_tibble(conf.matrix$table)
188
189 #library per costruire una confuzion matrix carina
190 library(cvms)
191
192 matrix<-table(tibble("target"=act,
193                      "prediction"=prd))
194
195 plot_confusion_matrix(as_tibble(matrix),
196                       target_col = "target",
197                       prediction_col = "prediction",
198                       counts_col = "n",
199                       add_normalized = FALSE,
200                       add_col_percentages = FALSE,
201                       add_row_percentages = FALSE,
202                       palette = "Greens")
203
204 #(ADI)
205 adjustedRandIndex (act,prd)
206
207 #visualizzo grafico con unità misclassificate e livello incertezza
208 library(mclust)
209 mis<-classError(act,prd)
210
211 #creo un dataset contenente variabili misclassified
212 error_data<-data[test_lab[mis$misclassified],]
213 error_data<-cbind(error_data,prd[mis$misclassified])
214 #ppca graphic of predicted variables
215 data_pred<-dada[test_lab,]
216
217 pca<-princomp(data_pred,cor=T)
218 prd = factor ( pred.values@partition , levels = c ("1","2") ,
219               labels = c (" B ", "M"))
220
221 miss_class<-mis$misclassified
222 data_pred$diagnosis<-as.character(prd)
223 data_pred$diagnosis[miss_class]<-"Errorclass"
224 data_pred$diagnosis<-as.factor(data_pred$diagnosis)

```

```

226 #grafico tramite factoextra che visualizza le prime due PCA
227 fviz_pca_ind(pca, col.ind=data_pred$diagnosis,
228             geom = c("point"),alpha=0.8,
229             pointsize=Uncertainty,
230             legend.title = "Diagnosis",
231             repel = TRUE,
232             title=c("2Dim PCA-EDDA Classifier"))+
233     scale_color_manual(values = c("palegreen2","gray","sandybrown"))
234
235 #-----
236 #FINITE MIXTURE of REGRESSION MODEL (GAUSSIAN)
237 #-----
238
239 #rileggo i dati
240 data<-read.csv("pancreatic_cancer.csv")
241
242 #data manipulation
243 data$sex<-ifelse(data$sex=="M",1,2)
244 #in questo caso consideriamo k =3
245
246 #PPCA imputation method
247 library(missMDA)
248 #ppca funz solo con num variabili
249 #ipotizza normalità delle variabili
250
251 data_num<-data[,-c(1,2,3,6,7,8)]
252 nb<-estim_ncpPCA(data_num,scale=T)
253 comp<-imputePCA(data_num,ncp = 5,scale=T)
254
255 data.pca<-abs(as.data.frame(comp$completeObs))
256 #view(data.pca)
257
258 data_num2<-data[, -c(1,2,3,7,8,9,14)]
259
260 #data visualization (non presente nel report)
261 #la variabilità aumenta man mano che la diagnosi peggiora
262 ggplot(data=data_num2, mapping=aes(x=creatinine,y=LYVE1))+
263     geom_point(aes(colour=diagnosis))+
264     geom_smooth(method="lm", se=F, size=1.5)+
265     theme_light()
266
267 library(mixtools)
268 library(flexmix)
269 # flexmix(formula, data, k = NULL, cluster = NULL,
270 #         model = NULL, concomitant = NULL,.....)
271
272 #aggiungo la variabile diagnosis
273 #nostra variabile risposta della regressione
274 datapi<-cbind(data[c(6)],data.pca)
275
276 #rimuovo outliers
277 data_ben=datapi[datapi$diagnosis==c(2),]
278
279 (outliers_tf<-boxplot.stats(data_ben$TFF1)$out)
280 (outliers_pla<-boxplot.stats(data_ben$plasma_CA19_9)$out)
281

```

```

282 data_fin<-datapi[-which(data_ben$TFF1 %in% outliers_tf),]
283 data_fin2<-data_fin[-which(data_ben$plasma_CA19_9 %in% outliers_pla),]
284 data_fin3<-data_fin2[-which(data_ben$plasma_CA19_9 %in% outliers_pla),]
285
286 datapi<-data_fin3
287 #str(datapi)
288 #algoritmo controllo E-M
289 itermax <- 1000
290
291 #entropia
292 eicval=1
293 eics<-matrix(nrow=itermax,ncol=2)
294 #bic
295 bicval <- Inf
296 bics<-matrix(nrow=itermax,ncol=2)
297 #n.b: da considerare abs(), potrebbe assumere valori negativi
298
299 #meccanismo di controllo per EM
300 #(tempistiche: circa 5 min)
301
302 #per ridurre il tempo ho confrontato solamente entropia
303 #tramite un metodo di MC, non molto ufficiale,
304 #se EIC si discosta molto dalla soglia 0.6 la classificazione
305 #inizia ad accorpare le due componenti M-B
306
307 for (iter in 1: itermax)
308 {
309   fit <- flexmix(diagnosis~.,data=datapi,k=3)
310   eics[iter,]<-c(iter,EIC(fit))
311   if (EIC(fit)<eicval)
312   {
313     eicval<-EIC(fit)
314     bestfit <- fit
315   }
316 }
317
318 #results:
319 summary(bestfit)
320 plot(bestfit)
321 BIC(bestfit)
322 EIC(bestfit)
323 ICL(bestfit) #ICL=entropia+BIC
324 KLdiv(bestfit) #distance Kullback Leiber
325
326 fit=bestfit
327
328 lab<-as.factor(datapi[,1])
329 pred<-as.factor(fit@cluster)
330
331 #riaggiusto i livelli etichette
332 levels(pred)<-c(2,3,1)
333 #pred
334 #lab
335 levels(pred)<-c("B","M","C")
336 levels(lab)<-c("C","B","M")
337

```

```

338 conf.matrix = confusionMatrix ( data = lab, reference = pred)
339 conf.matrix
340 #visualizza graficamente confuzion matrix
341 tab_conf<-as_tibble(conf.matrix$table)
342
343 matrix<-table(tibble("target"=lab,
344                     "prediction"=pred))
345
346 plot_confusion_matrix(as_tibble(matrix),
347                       target_col = "target",
348                       prediction_col = "prediction",
349                       counts_col = "n",
350                       add_normalized = FALSE,
351                       add_col_percentages = FALSE,
352                       add_row_percentages = FALSE,
353                       palette = "Oranges")
354
355 adjustedRandIndex(lab,pred)
356 fit@components
357 #varianza delle components
358 #come si notava dal grafico iniziale
359 #sono più elevate man mano che la diagnosi è più grave
360 (sigma<-parameters(fit)[10,])
361
362 #analisi della regressione
363 summary(refit(fit))
364
365 fit@components
366 fit@df #gradi libertà
367 #parametri
368 par<-parameters(bestffit)
369
370 #creo dataset con dati predicted
371 data_pred<-cbind(datapi[, -1],pred)
372 levels(pred)<-c("B","M","C")
373
374 #visualizzo funzione di densità per etichette previste
375 #considero LIVE1 in quantouna delle variabili più significative
376 ggplot(data_pred,aes(x=LYVE1,color=pred))+
377   geom_density(lwd=1)+
378   xlim(-10,30)+
379   scale_color_manual(values=c("green3", "sandybrown", "steelblue"))+
380   labs ( title = " FMM density - Pred")+
381   theme_bw()
382
383
384 datapi$diagnosis<-as.factor(datapi$diagnosis)
385 levels(datapi$diagnosis)<-c("C","B","M")
386 #visualizzo distribuzione dei gruppi reali
387 ggplot(datapi,aes(x=LYVE1,color=diagnosis))+
388   geom_density(lwd=1)+
389   xlim(-10,30)+
390   theme_bw()+
391   scale_color_manual(values = c("steelblue","palegreen2","sandybrown"))+
392   labs ( title = " FMM density - Real")
393

```