

MAE 5776

ANÁLISE MULTIVARIADA

Júlia M Pavan Soler
pavan@ime.usp.br

1º Sem/2020 - IME

Análise Multivariada

😊 Já vimos

- Análise Descritiva Multivariada
- Elipsóides de Dispersão e de Confiança, MANOVA
- Metodologias Clássicas de redução de dimensionalidade: $Y_{n \times p}$; $\mathbb{R}^p \rightarrow \mathbb{R}^m$
 - ✓ ACP ($m \leq \min(n, p)$), ACoP ($m \leq \min(n, p)$), AC ($m \leq \min(I-1, J-1)$), AF ($m \leq \min(n, p)$)
 - ✓ AG
 - ✓ AD ($m \leq \min(n, p, G-1)$), ACC ($m \leq \min(n, p, q)$)
 - ✓ Soluções Duais ($\mathbb{R}^{n \times p}$, $\mathbb{R}^{p \times p}$, $\mathbb{R}^{n \times n}$), Representações Biplot
 - ✓ PCR (Regressão via CP), PLS (Regressão via MQ Parciais)
- Integração de Bancos de Dados
- Diagrama de Caminhos ou Grafos

Var. quantitativas
 $n > p$
Obs iid



Caso: $n \ll p$ (*big p*): Soluções Regularizadas e Penalizadas

⇒ Componentes Principais
⇒ Análise Discriminante
⇒ Correlação Canônica

Métodos via Fatoração de Matrizes
e via Modelos de Regressão!

O Problema $n \ll p$

Big data
Big-p

- Dados "gasoline": $[Y_{60 \times 1} \quad X_{60 \times 401}]$

- Dados dos 3 Experts:

$$[Y1_{6 \times 3} \quad Y2_{6 \times 4} \quad Y3_{6 \times 3}] = Y_{6 \times 10}$$

Wine	Oak-type	Expert 1			Expert 2				Expert 3		
		Fruity	Woody	Coffee	Red fruit	Roasted	Vanillin	Woody	Fruity	Butter	Woody
1	1	1	6	7	2	5	7	6	3	6	7
2	2	5	3	2	4	4	4	2	4	4	3
3	2	6	1	1	5	2	1	1	7	1	1
4	2	7	1	2	7	2	1	2	2	2	2
5	1	2	5	4	3	5	6	5	2	6	6
6	1	3	4	4	3	5	4	5	1	7	5

- Dados dos Ratos Congenitos: $Y_{50 \times 35.129}$

- Dados dos Transcriptomas: $Y_{189 \times 22.215}$

```
library(devtools)
install_github("genomicsclass/tissuesGeneExpression")
library(tissuesGeneExpression)
data(tissuesGeneExpression)
dim(e) ## e contains the expression data
## [1] 22215 189
```

Redução de Dimensionalidade: $n > p \Rightarrow n \ll p$ Big data Big-p

X

1	nxp
2	
...	
n	

Componentes Principais e Coordenadas Principais:

$$\Sigma_{p \times p} = VDV' ; \quad F_k = XV_k$$

$$X_{n \times p} = UD^{1/2}V'; \quad F_k = U_k d_k^{1/2} \left\{ \max_a \frac{a' \Sigma a}{a' a}, \quad a' a = 1 \right.$$

$d_k > 0? \text{ (min}(n, p))$
 $V_{jk} = 0?$

X

Y

1	nxp
2	
...	
n	

1	Y
0	
...	
0 1	

Análise Discriminante (Linear de Fisher):

$$\Sigma_{p \times p} = \Sigma_B + \Sigma_W \Rightarrow \max_a \frac{a' \Sigma_B a}{a' \Sigma_W a}, \quad a' \Sigma_W a = 1$$

Σ_W inversível?

Y

X

1	nxp
2	
...	
n	

1	nxq
2	
...	
n	

Correlação Canônica:

$$\Sigma = \begin{pmatrix} \Sigma_{YY} & \Sigma_{YX} \\ \Sigma_{XY} & \Sigma_{XX} \end{pmatrix}$$

$$\left\{ \begin{array}{l} F_Y = a'Y; \quad \max_a \frac{a' \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} a}{a' \Sigma_{11} a}, \quad a' \Sigma_{11} a = 1 \\ F_X = b'X; \quad \max_b \frac{b' \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} b}{b' \Sigma_{22} b}, \quad b' \Sigma_{22} b = 1 \end{array} \right.$$

Σ_{11} inversível?
 Σ_{22} inversível?

PLS

$$\left[Cov(b'Y; a'X) \right]^2 = Var(b'Y) \left[Corr(b'Y; a'X) \right]^2 Var(a'X)$$

CCA

Componentes Principais Esparsos

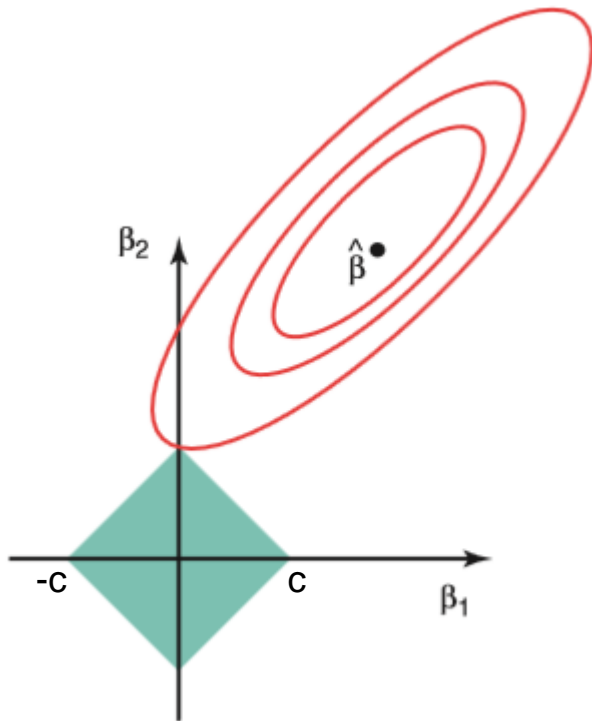
CP Penalizado (LASSO)

Penalização na forma de Lagrange:

$$\hat{\beta} = \arg \min_{\beta} \|F_k - Y\beta\|_2^2 + \lambda \|\beta\|_1$$

Penalização na forma de restrição:

$$\hat{\beta}_{2 \times 1}; \min_{\beta} \sum_{i=1}^n (F_{ik} - Y_i' \beta)^2 \text{ sujeito a } |\beta_1| + |\beta_2| \leq c$$



Solução mais esparsa: $\beta_1=0$

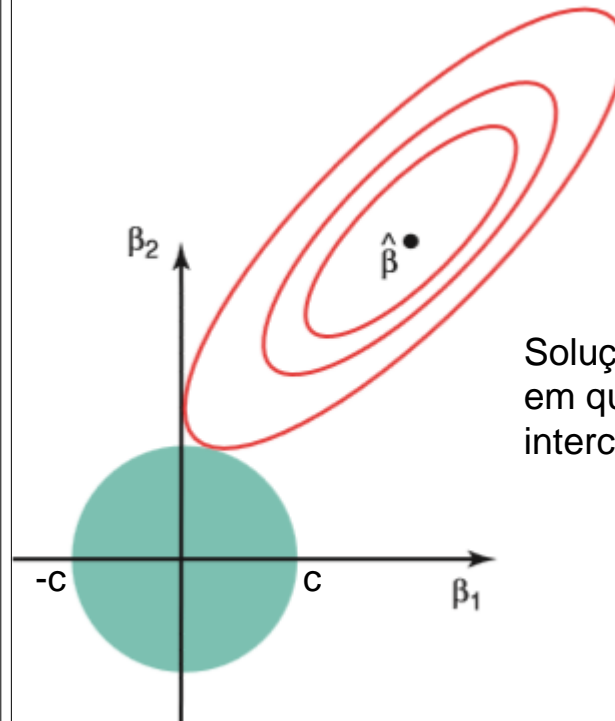
CP Regularizado (Ridge Regression)

Regularização na forma de Lagrange:

$$\hat{\beta} = \arg \min_{\beta} \|F_k - Y\beta\|_2^2 + \lambda \|\beta\|_2^2$$

Regularização na forma de restrição:

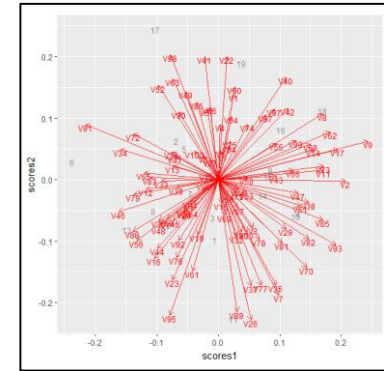
$$\hat{\beta}_{2 \times 1}; \min_{\beta} \sum_{i=1}^n (F_{ik} - Y_i' \beta)^2 \text{ sujeito a } \beta_1^2 + \beta_2^2 \leq c$$



Solução: primeiro ponto em que a elipse intercepta a restrição.

Solução menos esparsa: $\beta_1 \approx 0$

Componentes Principais em Espaços $n \ll p$ (*Big-p*)



Redução de dimensionalidade: Soluções regularizadas e penalizadas

$$Y_{n \times p} = \overset{\text{CoP}}{UD^{1/2}}V' \Rightarrow F_k = \overset{\text{CP}}{U_k d_k^{1/2}} = YV_k \quad n \ll p \Rightarrow \hat{F}_k = Y\hat{v}_k$$

$\uparrow \quad \uparrow$

Componente Principal Esparso (*Elastic Net*; Zou et al., 2006))

$$\hat{\beta} = \arg \min_{\beta} \left\{ \|F_k - Y\beta\|_2^2 + \lambda_1 \|\beta\|_2^2 + \lambda_2 \|\beta\|_1 \right\}; \quad \hat{v}_k = \frac{\hat{\beta}}{\|\hat{\beta}\|_2}; \quad \hat{F}_k = Y\hat{v}_k$$

$$\arg \min_{A, B} \left\{ \sum_{i=1}^n \left\| Y_{i_{p \times 1}} - A_{p \times m} B'_{m \times p} Y_i \right\|_2^2 + \lambda_1 \sum_{k=1}^m \|\beta_k\|_2^2 + \sum_{k=1}^m \lambda_{2k} \|\beta_k\|_1 \right\}; \quad A'A = I_m;$$

$$B_{p \times m} = (\beta_1, \dots, \beta_m)$$

Componente Principal Esparso via *svd* (Witten et al., 2009)

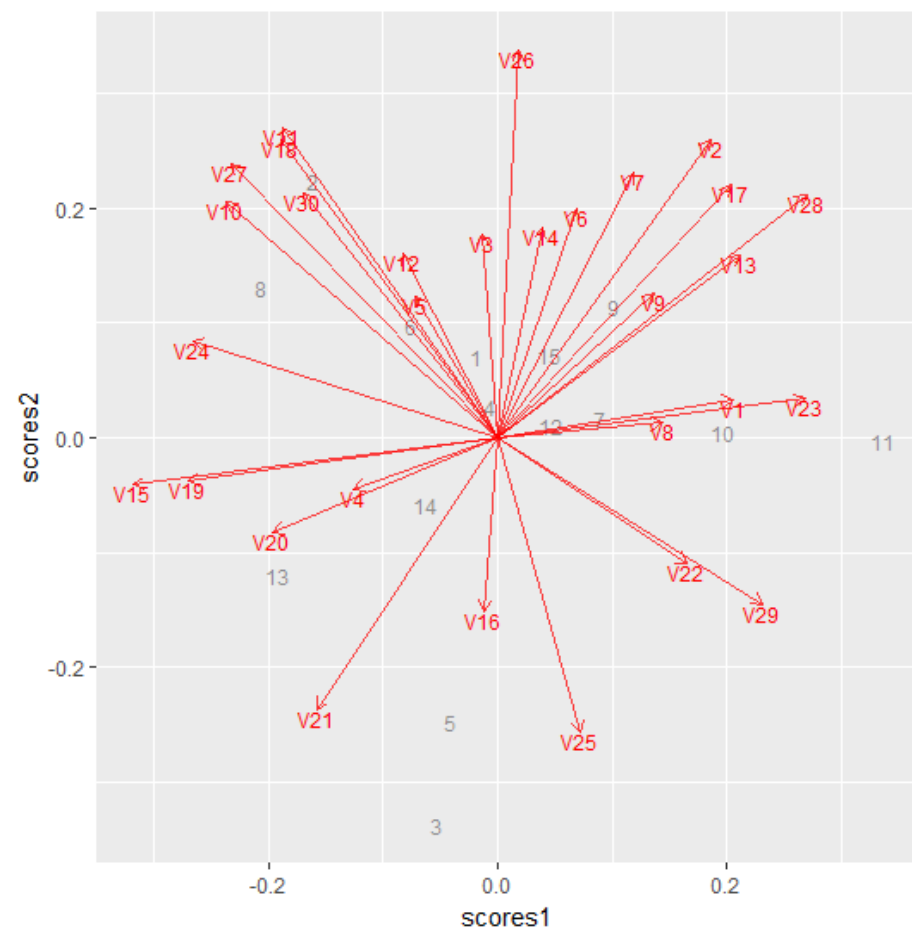
$$\max_{U_k, V_k} U_k' Y V_k; \quad \begin{cases} \|U_k\|_2^2 \leq 1 \\ \|V_k\|_2^2 \leq 1, \quad \|V_k\|_1 \leq c_1 \end{cases} \quad \tilde{F}_k = Y\tilde{V}_k$$

Componentes Principais – $n \ll p$

Biplot: $n=15$ $p=30$

CoP: Coordenadas Principais

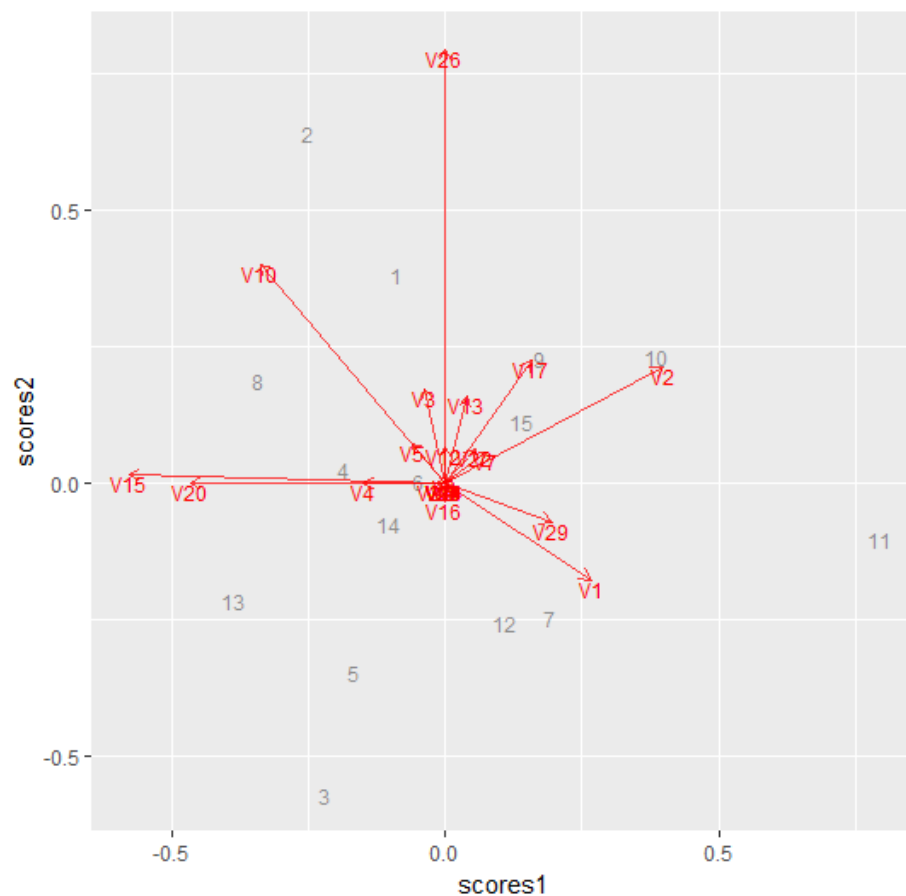
R-prcomp (suporta $n < p$)



Biplot: $n=15$, $p=30$

sCP: CP Esparso

R-SPCA do pacote ElasticNet



Big Data – $n \ll p$

Dados: Breast.TCGA do Bioconductor do R

```
breast.TCGA$data.train$mRNA
```

```
[1] 150 200 ← Big-p
```

```
breast.TCGA$data.train$miRNA
```

```
[1] 150 184 ← Big-p
```

```
breast.TCGA$data.train$proteomics
```

```
[1] 150 142 ←  $n > p$ 
```

```
breast.TCGA$data.train$subtype
```

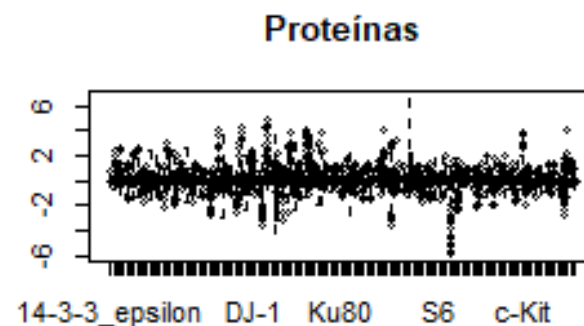
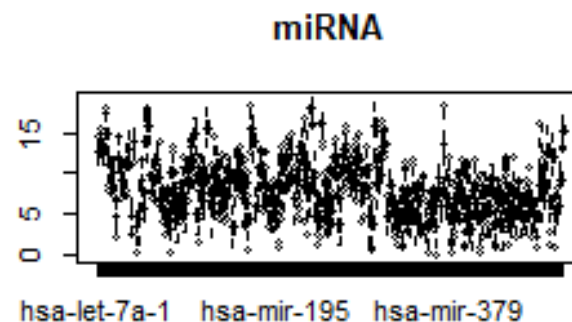
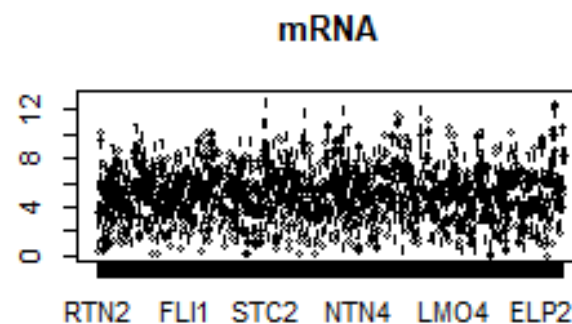
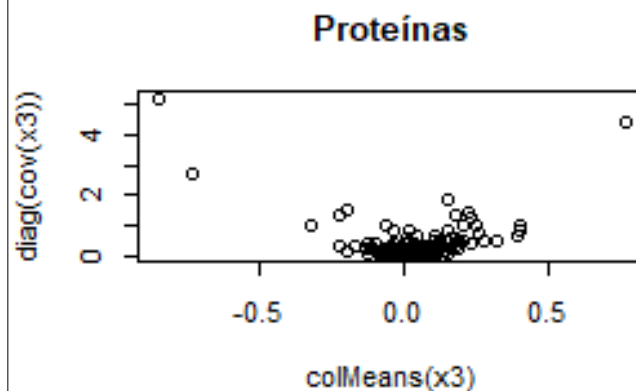
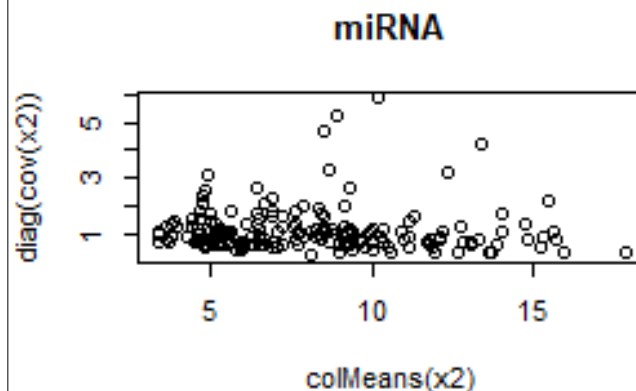
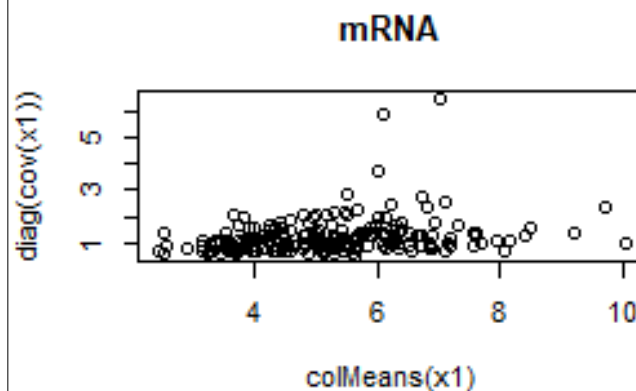
Basal	Her2	LumA
45	30	75

Dados
breast.TCGA

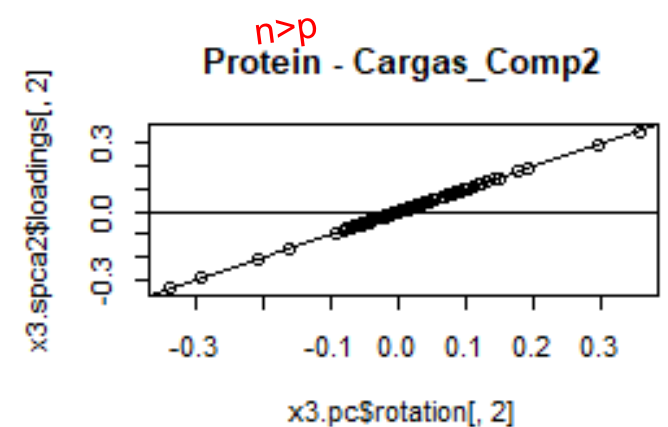
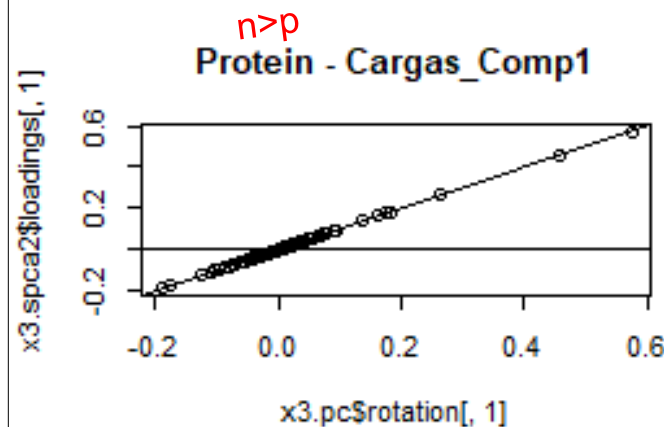
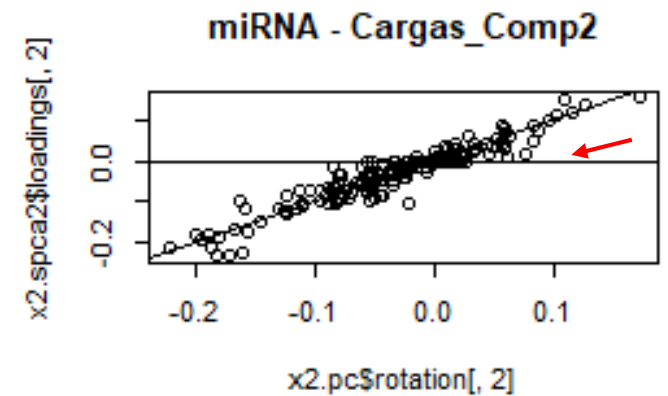
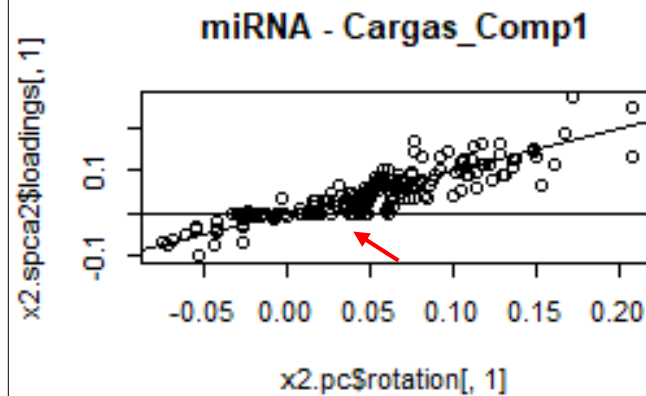
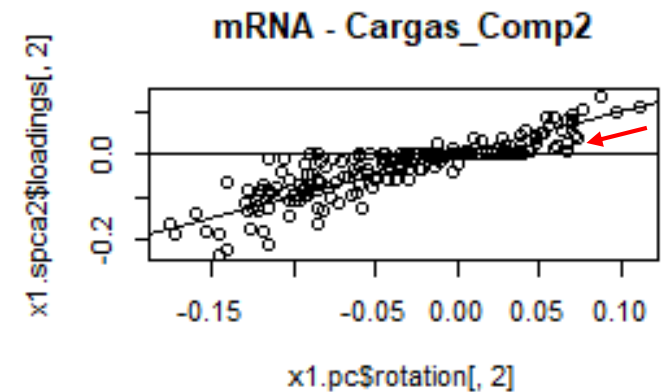
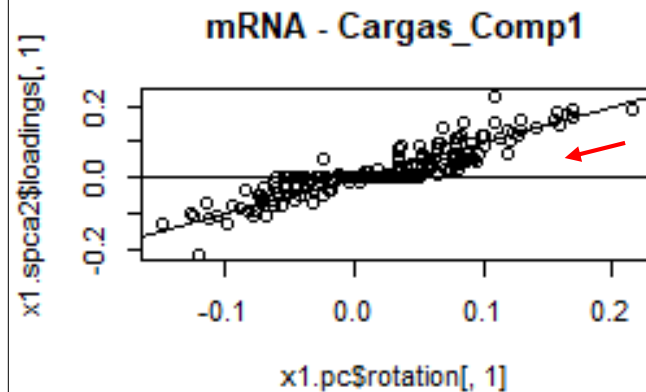
Dispersão:
média x var

Box-plots

Avaliação da
escala:
superdispersão,
obs atípicas



Componentes
Principais
(CoP: prcomp)

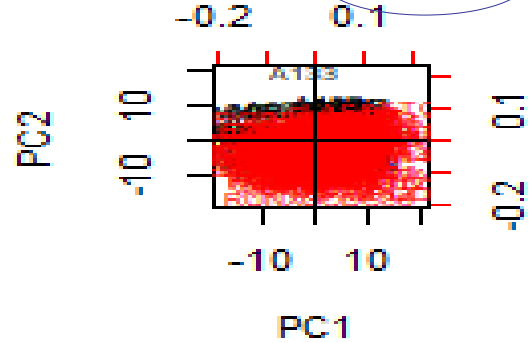


Comparação das
Cargas

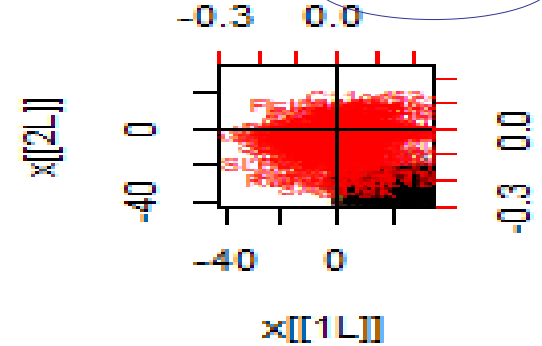
Componentes
Principais Esparsos
(sCP:elasticNet)

Biplots

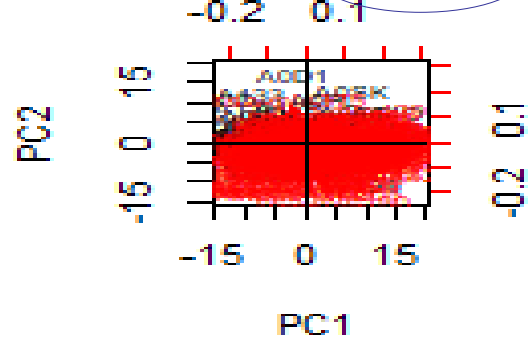
m.RNA - PC 34%



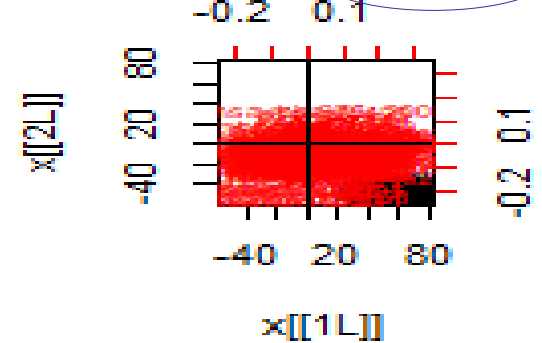
mRNA - sPC 12%



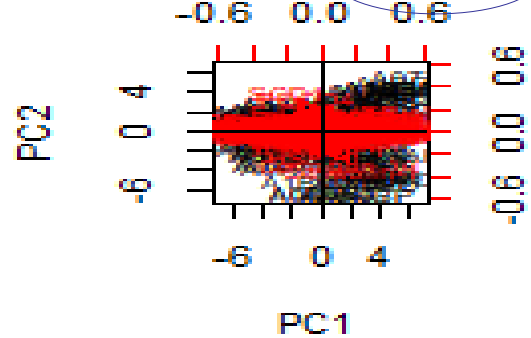
miRNA - PC 32%



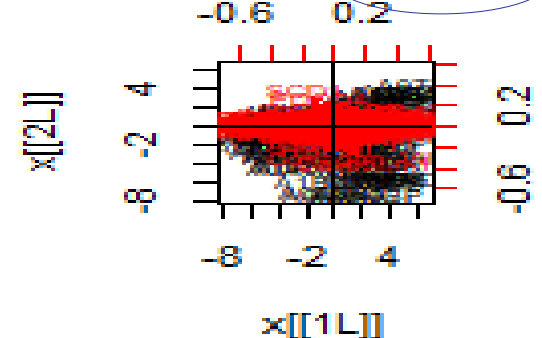
miRNA - sPC 28%



Protein - PC 39%



Protein - sPC 39%



- Indivíduos - Variáveis

Análise Discriminante Esparsa – $n \ll p$

Amostra do grupo g $Y_{i \times p} | \tau_g \stackrel{iid}{\sim} (\mu_g; \Sigma_g) \Rightarrow \hat{F}_i = l' Y_i \stackrel{iid}{\sim} (l' \mu_g; l' \Sigma l)$

Solução (linear) de Fisher: Suposição $\Rightarrow \Sigma_g = \Sigma$

Para $n > p$: S_W^{-1} ?

$$\max_l \frac{l' \sum_{g=1}^G n_g (\bar{Y}_g - \bar{Y})(\bar{Y}_g - \bar{Y})' l}{l' S_W l} = \max_l \frac{l' S_B l}{l' S_W l} \left\{ \begin{array}{l} \max_l l' S_B l ; \quad l_k' S_W l_k = 1, \\ l_j' S_W l_k = 0, j \neq k \\ k = 1, \dots, m \leq \min(n, p, G-1) \end{array} \right.$$

Tabela de MANOVA

$P_G = G(G'G)^{-1}G'$ $G_{n \times G}$: Matriz de incidência de grupos

F.V.	g.l.	Matriz de SQPC
Trat	G-1	$H_{p \times p} = \sum_{g=1}^G n_g (\bar{y}_g - \bar{y})(\bar{y}_g - \bar{y})' = S_B = Y'P_G Y$
Resíduo	n-G	$E_{p \times p} = \sum_{g=1}^G \sum_{i=1}^{n_g} (y_{gi} - \bar{y}_g)(y_{gi} - \bar{y}_g)' = S_W = Y'(I - P_G)Y$
TOTAL	n-1	$H + E = \sum_{g=1}^G \sum_{i=1}^{n_g} (y_{gi} - \bar{y})(y_{gi} - \bar{y})' = Y'Y$

Análise Discriminante Esparsa – $n \ll p$

$$\max_{\beta_k} \frac{\beta_k' S_B \beta_k}{\beta_k' S_W \beta_k}; \quad \left| S_W^{-1/2} S_B S_W^{-1/2} - d I_p \right| = 0 \quad \overset{n \ll p}{\Rightarrow} \quad S_W \text{ não é p.d (?)}$$

Tornar os autovalores positivos

$$\max_{\beta_k} \frac{\beta_k' S_B \beta_k}{\beta_k' (S_W + \Omega) \beta_k}; \quad \left| (S_W + \Omega) - d I_p \right| = 0, \quad d > 0$$

$$\Rightarrow \max_{\tilde{\beta}_k} \tilde{\beta}_k' \tilde{S}_B \tilde{\beta}_k \begin{cases} \tilde{S}_B = (S_W + \Omega)^{-1/2} S_B (S_W + \Omega)^{-1/2} \\ \beta_k' (S_W + \Omega) \beta_k = 1; \quad \beta_k' (S_W + \Omega) \beta_j = 0; \end{cases}$$

Hastie et al., (1995); Clemmensen et al., (2011, 2016)

Solução regularizada para corrigir o posto de S_W

Mas, como obter Ω ?
E os β 's, são esparsos?

Solução Regularizada e Penalizada: implementada em *sda* do pacote `sparseLDA_R`

Problema de predição de grupos

$$\hat{\beta}_k \text{ }_{p \times 1}; \quad \min_{\beta_k, \theta_k} \left\{ \left\| G_{n \times G} \theta_k - Y_{n \times p} \beta_k \right\|_2^2 + \lambda_1 \beta_k' \Omega \beta_k + \lambda_2 \left\| \beta_k \right\|_1 \right\}; \quad \theta_k' G' G \theta_k = 1, \quad \theta_k' G' G \theta_j = 0,$$

Vetor de pesos dos grupos

Para $\Omega = I_m$: usar algoritmo do *ElasticNet*

Análise Discriminante Esparsa – $n \ll p$

■ Big-Data: $n \ll p$

`sparseLDA`: Dados “penicilliumYES”, espécies de fungos (indistinguíveis) avaliadas em variáveis de imagem

$n=36$ $p=3.754$

$G=3$: "*P. Melanoconidium*", "*P. Polonicum*", "*P. Venetum*"

	$Y\hat{\beta}_1$	$Y\hat{\beta}_2$
[1,]	-3.113028	-3.4122173
[2,]	-3.142295	-3.8733571
[3,]	-2.988152	-1.4446112
[4,]	-2.995266	-1.5567002
[5,]	-2.995715	-1.5637771
[6,]	-3.063970	-2.6392373
[7,]	-2.996005	-1.5683428
[8,]	-2.998670	-1.6103476
[9,]	5.059007	0.9703178
[10,]	5.803007	0.7348022
[11,]	2.720020	0.5259276
[12,]	4.546290	-0.6951151
[13,]	6.415934	-0.6321676
[14,]	5.898611	0.1021021
[15,]	5.472054	-1.9047044
[16,]	8.543129	-0.2298407
[17,]	-2.687269	3.2962259
[18,]	-2.666652	3.6210784
[19,]	-2.840601	0.8802641
[20,]	-2.848656	0.7533365
[21,]	-2.733009	2.5755237
[22,]	-2.262692	2.7361895
[23,]	-2.779546	1.8422657
[24,]	-1.346527	3.0923848

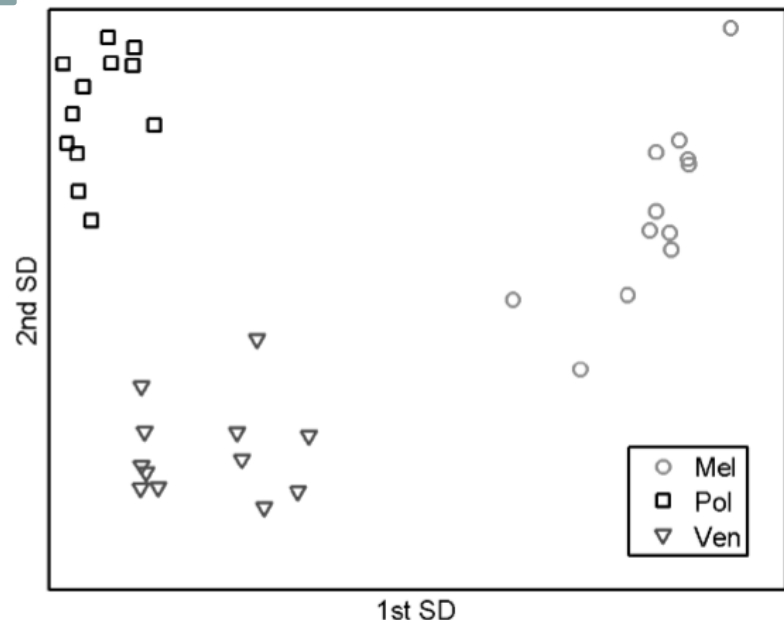
...

	$\hat{\theta}_1$	$\hat{\theta}_2$
[1,]	0.7357055	1.20778203
[2,]	-1.4138227	0.03324864
[3,]	0.6781172	-1.24103066

Matriz de pesos
dos grupos

Escores da funções
discriminantes

Representação das espécies
nas variáveis discriminantes



Correlação Canônica Esparsa – $n \ll p$

$$Y_{i(p+q) \times 1} = \begin{pmatrix} Y_{1i} \\ Y_{2i} \end{pmatrix} \stackrel{iid}{\sim} \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}; \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right) \Rightarrow \begin{cases} a_k' Y_1 \\ b_k' Y_2 \end{cases} \max_{a_k, b_k} \rho(a_k' Y_1, b_k' Y_2)$$

$$\Rightarrow \frac{a' \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} a}{a' \Sigma_{11} a} \Rightarrow \frac{b' \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} b}{b' \Sigma_{22} b}$$

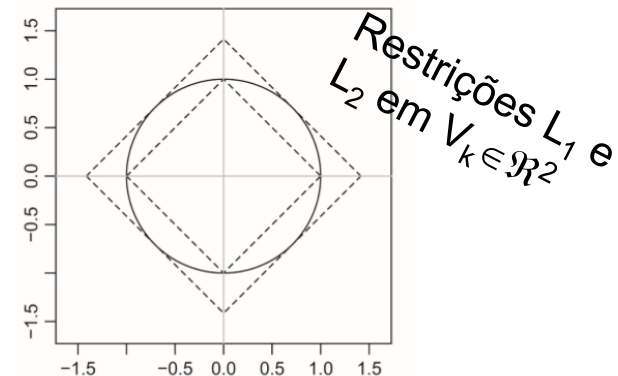
Problemas na otimização quando $n \ll p$ $n \ll q$

Inicialmente, considere o problema de obter uma solução penalizada para a decomposição em valores singulares (svd) de matrizes:

PMD: Penalized Matrix Decomposition (Witten, Tibshirani, Hastie, 2009, 2015)


$$Y_{n \times p} = U D V', \quad U' U = I_n, V' V = I_p$$

$$\min_{U_k, V_k, d_k} \left\| Y - U_k V_k' d_k \right\|_2^2; \quad \begin{cases} \|U_k\|_2^2 \leq 1, \|U_k\|_1 \leq c_1 \\ \|V_k\|_2^2 \leq 1, \|V_k\|_1 \leq c_2 \end{cases}$$



Decomposição de Matrizes - Penalização

PMD: Penalized Matrix Decomposition

$$\min_{U_k, V_k, d_k} \left\| Y - U_k V_k' d_k \right\|_2^2 ; \quad \begin{cases} \|U_k\|_2^2 \leq 1, \|U_k\|_1 \leq c_1 \\ \|V_k\|_2^2 \leq 1, \|V_k\|_1 \leq c_2 \end{cases}$$


$$\frac{1}{2} \left\| Y - U D V' \right\|_2^2 = \frac{1}{2} \|Y\|_2^2 - \sum_{k=1}^m U_k' Y V_k d_k + \frac{1}{2} \sum_{k=1}^m d_k^2$$

$$\max_{U_k, V_k} U_k' Y V_k ; \quad \begin{cases} \|U_k\|_2^2 \leq 1, \|U_k\|_1 \leq c_1 \\ \|V_k\|_2^2 \leq 1, \|V_k\|_1 \leq c_2 \end{cases} \quad \text{PMD}(L_1, L_1)$$

Diferentes algoritmos são propostos para solução deste problema de maximização (Witten et al., 2009, 20015)

Bilinear em U e V \Rightarrow U fixo e obter V, V fixo e obter U

$$V_k \text{ fixo} : \max_{U_k} U_k' Y V_k ; \quad \|U_k\|_2^2 \leq 1, \|U_k\|_1 \leq c_1, 1 \leq c_1 \leq \sqrt{n}$$

$$U_k \text{ fixo} : \max_{V_k} U_k' Y V_k ; \quad \|V_k\|_2^2 \leq 1, \|V_k\|_1 \leq c_2, 1 \leq c_2 \leq \sqrt{p}$$

Decomposição de Matrizes – Penalização Aplicação

■ CCA – Esparso: PMD(L_1, L_1)

$$\max_{a_k, b_k} a_k' \overset{\text{covariância}}{Y_1' Y_2} b_k ; \quad \begin{cases} a_k' Y_1' Y_1 a_k \leq 1, \quad \|a_k\|_1 \leq c_1 \\ b_k' Y_2' Y_2 b_k \leq 1, \quad \|b_k\|_1 \leq c_2 \end{cases}$$

Algoritmo proposto: **CCA-P Diagonal**. Assume que para dados em alta dimensão a matriz de covariância diagonal pode ser adotada (Dudoit et al. 2001; Tibshirani et al., 2003)

$$a_k' Y_1' Y_1 a_k \overset{Y_1' Y_1 = I_p}{=} a_k' a_k \leq 1, \quad b_k' Y_2' Y_2 b_k \overset{Y_2' Y_2 = I_q}{=} b_k' b_k \leq 1$$

■ PCA – Esparso: PMD(\cdot, L_1) (Witten et al., 2009, 2015)

$$\max_{V_k} V_k' Y' Y V_k ; \quad \|V_k\|_2^2 \leq 1, \quad \|V_k\|_1 \leq c_2$$

Correlação Canônica Esparsa

- Pacote PMA-R, função CCA_R:

Dados gerados: $n=100$ $p=500$ $q=1000$: $\begin{bmatrix} Y_{100 \times 500} & Y_{200 \times 1000} \end{bmatrix}$

$$U_{n \times 1} = Y_1 a \quad \Rightarrow \quad n(a_j \neq 0) = 338;$$

$$V_{n \times 1} = Y_2 b \quad \Rightarrow \quad n(b_j \neq 0) = 687; \quad c_1 = c_2 = 0,5667$$

$$\rho_c = 0,9735$$

```
Set.seed(19)
Set.seed(90)
Num non-zeros u's: 338
Num non-zeros v's: 687
Type of x: standard
Type of z: standard
Penalty for x: L1 bound is 0.5666667
Penalty for z: L1 bound is 0.5666667
Cor(Xu, Zv): 0.9735552
```