MAE 5776

ANÁLISE MULTIVARIADA

Júlia M Pavan Soler

pavan@ime.usp.br

Análise Multivariada

$$Y_{n\times p} = (Y_{ij}) \in \Re^{n\times p}$$



- Estatísticas Descritivas Multivariadas
- ✓ Distribuição Normal Multivariada, Distribuições Amostrais
- Regiões de Confiança, Testes Multivariados, MANOVA, IC Simultâneos, Correções para Múltiplos Testes
- ✓ Análises Multivariadas Clássicas (n>p, iid, ℜ^{nxp}): CP, CoP, AC, AF, AD, ACC

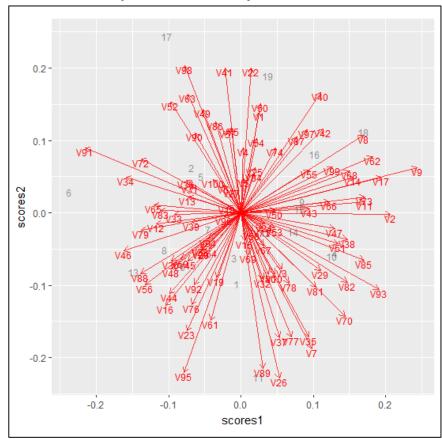
Análises Multivariadas em Dados de Alta Dimensão (Big-Data) Razão n/p?

- ✓ Análises Multivariadas em Big-p (n<<p): Soluções duais e Soluções regularizadas e esparsas (CP, AD, ACC)
- ⇒ Análises Multivariadas em Big-n (n >> p)

Dados Multivariados: $Y_{n \times p}$ Dimensionalidade dos Dados

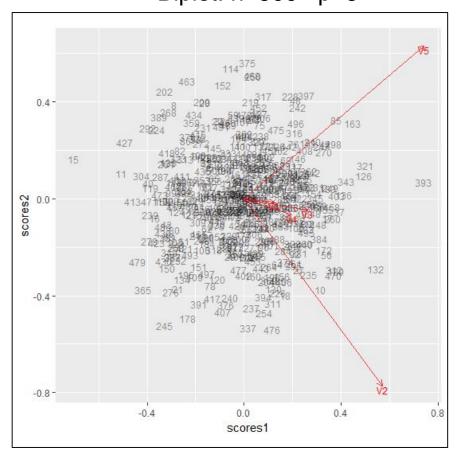
Big-p (n << p) "wide data"

Biplot: n=20 p=100



Soluções em espaços duais Soluções regularizadas e/ou penalizadas

Big-n (n >> p) "long data" Biplot: n=500 p=5



Sumarização e visualização ? (Black Screen Problem: *R_alpha blending*)

Não confundir Estrutura de Dados Dados: Medidas Repetidas com outras no formato "long" definições!

Dados: Medidas Repetidas no formato "wide"

Alternativa de análise: MANOVA > library(reshape2) > dat.wide <- dcast(dat.long, Subj +</pre>

	, aac.wiac (acase (aac. rong, sas)			
Grı	<pre>Grup + Staph~ > dat.wide</pre>		Time, value.var="		="02")	
> (
	Subj	Grup	Staph	6	12	18
1	1	P	1	1.48	2.81	3.56
2	2	P	0	1.04	2.07	2.81
3	3	P	1	1.48	2.52	3.41
4	4	P	0	1.04	1.93	2.89
21	21	P	1	1.50	2.85	3.12
22	22	Р	0	1.65	2.70	3.40
23	23	Р	1	1.80	2.15	3.90
24	24	Р	0	1.20	2.25	3.30
25	25	V	1	1.78	2.96	4.00
26	26	V	0	1.48	2.81	3.85
27	27	V	1	1.33	2.52	3.84
28	28	V	0	1.03	2.07	2.96
45	45	V	1	1.65	3.00	4.05
46	46	V	0	1.20	2.70	3.90
47	47	V	1	1.35	2.55	3.67
48	48	V	0	1.20	2.70	3.60

no formato "long" > library (MANOVA.RM) > dat.long<-o2cons Alternativa de									
> li	> library (MANOVA.RM) > dat.long<-o2cons > dat.long > dat.long O2 Staph Time Grup Subj								
> da	> dat.long<-o2cons				ling	***Se:	Mo. de		
> da	t.lon	ıg			.,69	ares r	L'odelos		
> da	t.lon			1		- 7	"Stos		
		Staph		Grup	Subj				
	1.48	1	6	Р	1				
	2.81		12	Р	1				
	3.56	1		Р	1				
	1.04	0	6	Р	2				
	2.07	0	12	Р	2				
	2.81	0	18	Р	2				
67	1.80	1	6	P	23				
	2.15	1	12	- P	23				
	3.90	1	18	P	23				
	1.20	0	6	P	24				
	2.25	0	12	Р	24				
	3.30	0	18	Р	24				
	1.78	1	6	V	25				
74	2.96	1	12	V	25				
75	4.00	1	18	V	25				
76	1.48	0	6	V	26				
77	2.81	0	12	V	26				
78	3.85	0	18	V	26				
130	1.35	1	6	V	47				
	2.55	1	12	V	47				
	3.67	1	18	V	47				
	1.20	0	6	V	48				
	2.70	0	12	V	48				
	3.60	0	18	V	48				

Dados Multivariados - O que é Big-Data ?

- Hair et al., 2005, 2018: Tamanho amostral ⇒ n>100, n > 5p
- Grau de correlação entre as variáveis ⇒ |r| > 30%

Análises Clássicas

- Fokoué, E. (2015): apresenta uma taxonomia para dados de "Alta Dimensão"
 - → Grandes Bancos de Dados: n>1000 ou p>50

Razão n/p!

	$\frac{n}{n} < 1$	$1 \leq \frac{n}{p} < 10$	$\frac{n}{p} \geq 10$	
	$\begin{array}{c} p \\ \textbf{Information} \\ \textbf{Poverty} \\ (n <\!\!<\!\!< p) \end{array}$	Information Scarcity	$\begin{array}{c} p - \\ \textbf{Information} \\ \textbf{Abundance} \\ (n \ggg p) \end{array}$	
n > 1000	Large p , Large n	Smaller p , Large n	Much smaller p , Large n	
<i>n</i> ≤ 1000	Large p , Smaller n	Smaller p , Smaller n	Much smaller p , Small n	

In this taxonomy, **A** *and* **D** *pose a lot of challenges.*

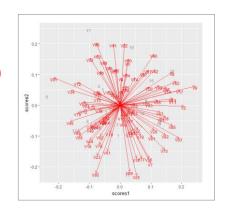
Já vimos!

Dados Multivariados – Big-p



Redução de dimensionalidade: Soluções duais, soluções regularizadas e penalizadas

$$Y_{n \times p} = U \Lambda^{1/2} V' \stackrel{n << p}{\Rightarrow} Z_j = U_j d_j^{1/2} \Rightarrow \hat{Z}_j = Y \hat{v}_j$$



Componente Principal Regularizado e Penalizado: (Elastic Net)

(Zou, Hastie and Tibshirani, 2006)

CoP (solução dual)

$$\hat{\beta} = \arg\min_{\beta} \left\{ \left\| Z_{j} - Y\beta \right\|_{2}^{2} + \left\| \lambda_{1} \left\| \beta \right\|_{2}^{2} \right\| + \left\| \lambda_{2} \left\| \beta \right\|_{1} \right\}; \quad \hat{v}_{j} = \frac{\hat{\beta}}{\left\| \hat{\beta} \right\|_{2}}; \quad \hat{Z}_{j} = Y \hat{v}_{j}$$

$$\arg\min_{\alpha,\beta} \sum_{i=1}^{n} \left\| Y_{i_{p\times 1}} - \alpha_{p\times 1} \beta_{1\times p}^{i} Y_{i} \right\|_{2}^{2} + \lambda \left\| \beta \right\|_{2}^{2}; \quad \left\| \alpha \right\|_{2}^{2} = 1 \quad \Rightarrow \hat{\beta} \propto \hat{v}_{1}$$

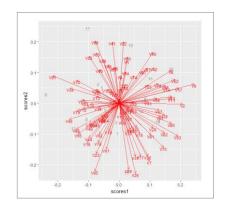
Já vimos!

Dados Multivariados – Big-p



Redução de dimensionalidade: Soluções duais, soluções regularizadas e penalizadas

$$Y_{n \times p} = U \Lambda^{1/2} V' \stackrel{n << p}{\Rightarrow} Z_j = U_j d_j^{1/2} \Rightarrow \hat{Z}_j = Y \hat{v}_j$$



Decomposição em valores singulares de matrizes

(Witten et al., 2009; Tibshirani et al., 2015)

$$\max_{U_{k},V_{k}} U_{k}^{\prime} Y_{n \times p} V_{k} ; \begin{cases} \left\| U_{k} \right\|_{2}^{2} \leq 1, \left\| U_{k} \right\|_{1} \leq c_{1} \\ \left\| V_{k} \right\|_{2}^{2} \leq 1, \left\| V_{k} \right\|_{1} \leq c_{2} \end{cases}$$
PMD(L₁,L₁)

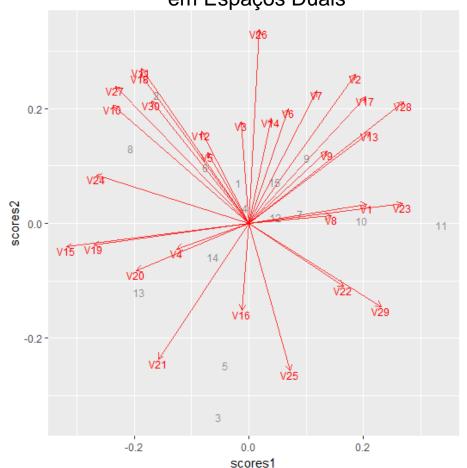
- $\Rightarrow \mathbf{CP}: \quad \max_{V_k} V_k'(Y'Y)_{p \times p} V_k; \quad \|V_k\|_2^2 \le 1, \quad \|V_k\|_1 \le c_2 \quad \quad \mathsf{PMD}(\cdot, \mathsf{L}_1)$
- \Rightarrow CoP: $\max_{U_k} U_k'(YY')_{n \times n} U_k; ||U_k||_2^2 \le 1, ||U_k||_1 \le c_1 \quad PMD(L_1, \cdot)$
- $\Rightarrow ACC: \max_{a_k, b_k} a_k' (Y_1' Y_2)_{p \times q} b_k ; \begin{cases} a_k' Y_1' Y_1 a_k \le 1, \ \|a_k\|_1 \le c_1 \\ b_k' Y_2' Y_2 b_k \le 1, \ \|b_k\|_1 \le c_2 \end{cases} PMD(L_1, L_1)$

Já vimos!

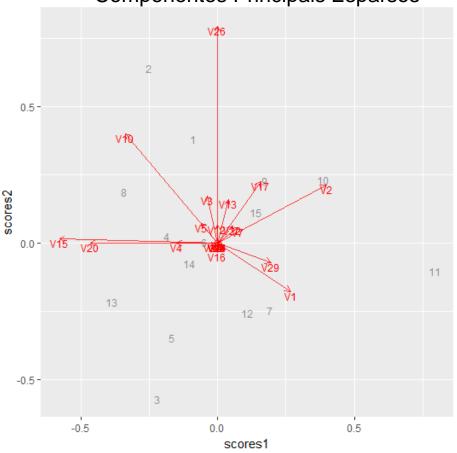
Componentes Principais – n<<p

Representação Biplot: n=15 p=30

R-prcomp: Componentes Principais em Espaços Duais



R-SPCA do pacote ElasticNet: Componentes Principais Esparsos



Dados Multivariados — Big-n

Big-n ⇒ erro padrão dos estimadores tendem a zero

$$\frac{s}{\sqrt{n}} \to 0$$
 \notin problema inferencial, somente análise descritiva de dados ?? Não é um consenso!

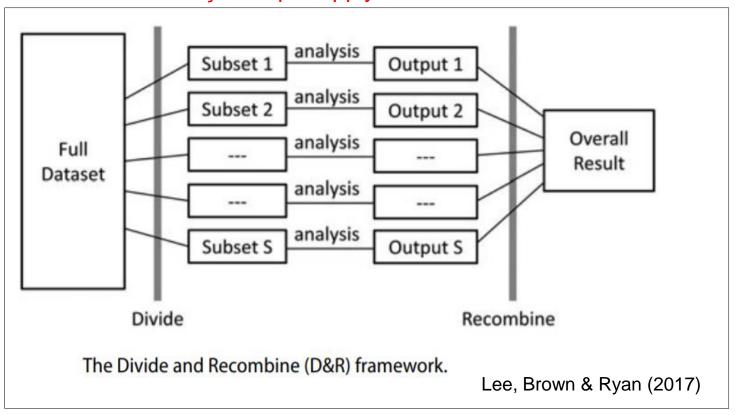
```
No R: read.table() ⇒ read.big.matrix() library(bigmemory) library(biganalytics)
```

- Análise de Dados em Big-n ⇒ Soluções Split-Apply-Combine ou Divide-Recombine
 - SAS: by
 - Google: MapReduce
 - Apache: Hadoop
 - Bibliotecas do R: biganalytics

```
partools (Software Alchemy-Matloff, 2016) foreach, doMC (multiple Cores)
```

Dados Multivariados – Big-n

Racional das Soluções Split-Apply-Combine ou Divide-Recombine:

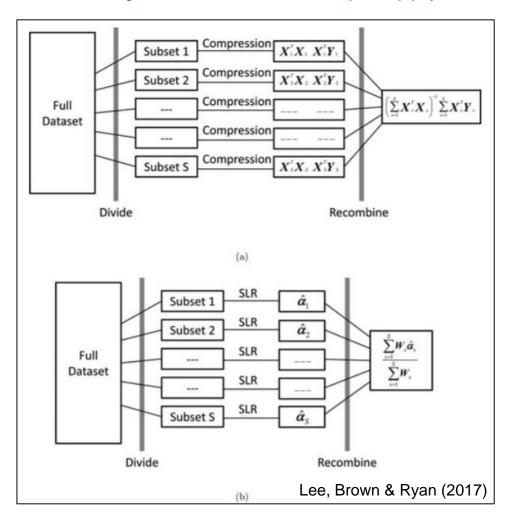


Soluções Paralelizadas (de cálculos estatísticos em amostras "iid"):

- Particionar os dados via Aleatorização (com repos.) ou condicional a uma variável
- Em cada sub-amostra calcular a estatística de interesse
- Obter a medida "agregada" das estatísticas nas subamostras

Dados Multivariados – Big-n

Soluções Paralelizadas Splt-Apply-Combine ou Divide-Recombine (D&R):



Propriedades assintóticas dos EMVS na família exponencial (Portnoy, 1988)

$$\hat{\alpha}_{EMVS} \stackrel{\frac{p^2}{n} \to 0; p, n \to \infty}{\sim} N(\alpha; v(\alpha))$$
Limite inferior de Kramer-Rao

$$n = \sum_{s=1}^{S} n_{s}; \quad n_{s} = \frac{n}{S}$$

$$\Rightarrow \bar{\alpha}_{n} = \frac{1}{S} \sum_{s=1}^{S} \hat{\alpha}_{s} = \frac{\sum_{s=1}^{S} W_{s} \hat{\alpha}_{s}}{\sum_{s=1}^{S} W_{s}}$$

$$W_{s} = X_{s}' X_{s}$$

$$\Rightarrow Cov(\bar{\alpha}_{n}) = \frac{1}{S} Cov(\hat{\alpha}_{s})$$
D&R horizontal em Ajustes lineares
$$Ajustes lineares$$

$$S \text{ fixo} \\ n \to \infty, n_{s} \to \infty$$

$$\bar{\alpha}_{n} \to \hat{\alpha}_{EMVS}$$

Para análises lineares (*lm*) tais soluções (D&R horizontal e via suficiência) são equivalentes. Para *glm* isto não ocorre!

(a) D&R via estatísticas resumo (estat. suficiente)

(b) D&R horizontal: método clássico

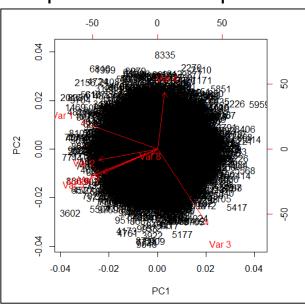
```
require (partools)
                                    Dados Multivariados – Big-n
cls <- makeCluster(2)</pre>
setclsinfo(cls)
                                                         Matloff (2015):
#Gerar dados (n>>p) para ajustar Modelo Linear
                                                         solução paralelizada
n <- 10000`
               n/p=5000
p <- 2
set.seed = 1099
tmp <- matrix(rnorm((p+1)*n),nrow=n)</pre>
u <- tmp[,1:p] # gerar valores "X"
# adicionar coluna de valores "Y"
u \leftarrow cbind(u, u % * p (1, p) - t tmp[, p+1])
colnames(u) = c("X1", "X2", "Y")
                                              > head(u)
#head(u)
                                                        Х1
                                                                 X2
# ajustar lm via solução paralelizada
                                              [1,] -0.2056215 -0.2229407 -1.712522
# D&R horizontal
                                              [2,] -0.5722810 -1.2393545 -1.965369
                                              [3,] 0.5898422 -0.4974834 1.028857
# (N. Matloff, 2015)
                                              [4,] 0.6045709 -0.9882614 -1.649882
# SA: Software Alchemy
                                              [5,] 0.4346162 0.8333716 1.488213
distribsplit (cls, "u")
                                              [6,] 0.1783079 -0.6069011 -1.326253
\#calm(cls, "u[,3] \sim u[,1]+u[,2]")
calm(cls, "u[,3] \sim u[,1]+u[,2]")$tht
 (Intercept) u[, 1] u[, 2]
-0.003110128 1.005448362 1.002561369
# check: resultados devem ser aproximadamente os mesmos
lm(u[,3] \sim u[,1]+u[,2])
 (Intercept) u[, 1]
                               u[, 2]
  -0.002909 1.005829
                               1.002436
```

Dados Multivariados – Big-n

```
require(partools)
cls <- makeCluster(2)
setclsinfo(cls)

#gerar dados multivariados
library(clusterGeneration)
n = 10000
p = 8
n/p
[1] 1250
media_pop = c(rdunif(p, 10))
cov_pop = round(genPositiveDefMat(dim=p)$Sigma,2)
dados_=-mvrnorm(n, media_pop, cov_pop)
biplot(prcomp(dados, scale=TRUE))</pre>
```

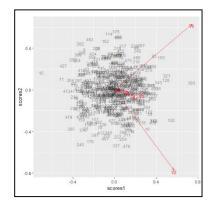
Biplot: Dados completos



```
#Rodar análise de CP via solução paralelizada (N. Matloff, 2015)
#SA: Sottware Alchemy
distribsplit(cls, "dados")
caprcomp(cls, 'dados, scale=TRUE', 8) $sdev
[1] 1.4707854 1.0525281 1.0046278 0.9884619 0.9051253 0.8954517
[7] 0.8450872 0.6382045
prcomp(dados, scale=TRUE) $sdev
[1] 1.4707662 1.0516048 1.0024852 0.9883842 0.9063294 0.8971400
[7] 0.8458021 0.6382948
```

Desvio padrão dos 8 primeiros CP são "iguais"

Dados Multivariados – *Big-n* Visualização



Visualização de Dados de Alta Dimensão – Big_n (Norman Matloff)

Dificuldade: "Borrão" (BSP)



Construção do Gráfico em Coordenadas Paralelas

Representação dos pontos (perfis) mais frequentes: Padrões TopFrequence

Dados Contínuos: Estimar a densidade (método dos *k*-vizinhos mais próximos) e Representar os padrões TopFrequence ou os Outliers

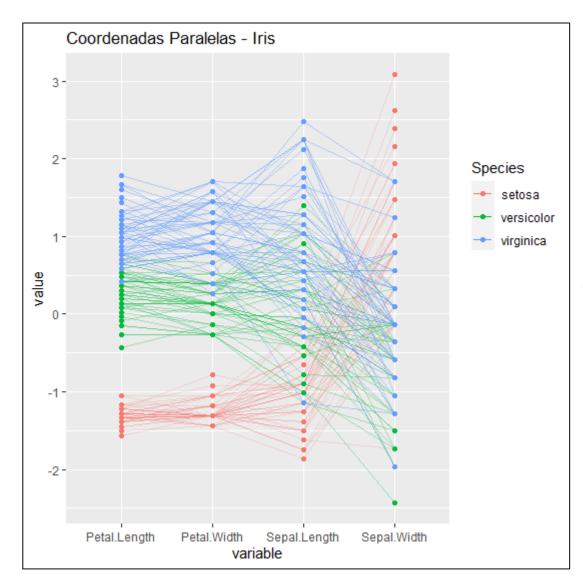
Dados Heterogêneos: Categorizar os dados contínuos e Representar os Padrões TopFrequence ou os Outliers

Gráfico em Coordenadas Paralelas

Dados - Iris

 $Y_{150x(4+1)}$

Valores: escore z

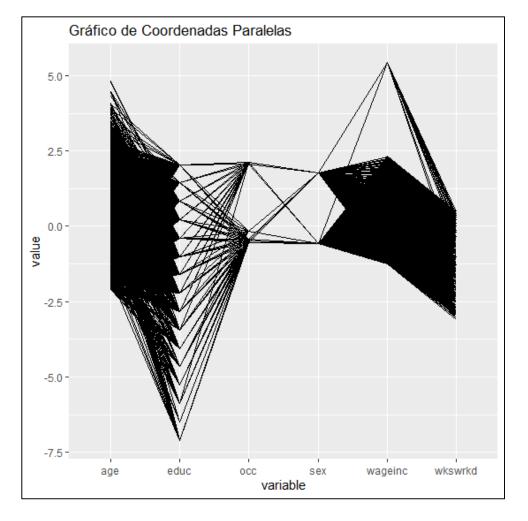


Boa visualização ("borrão" não aparece!)

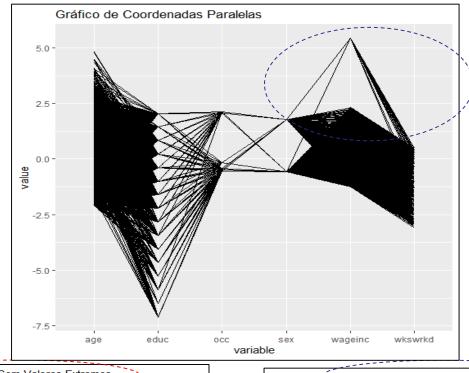
Big-n

Visualização

```
> data(prgeng)
> pe <- prgeng[,c(1,3,5,7:9)] # extrair 6 vars de interesse
> dim(pe)
[1] 20090
              6
> head(pe)
      age educ occ sex wageinc wkswrkd
1 50.30082
             13 102
                      2
                          75000
                                     52
                                     20 ...
2 41.10139
              9 101
                          12300
```



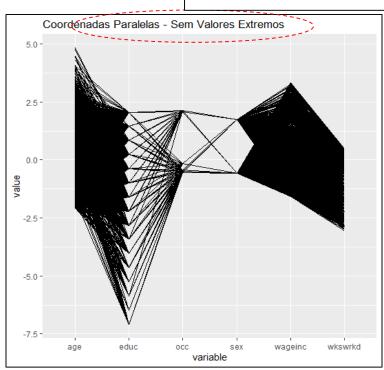
Problema: "borrão"

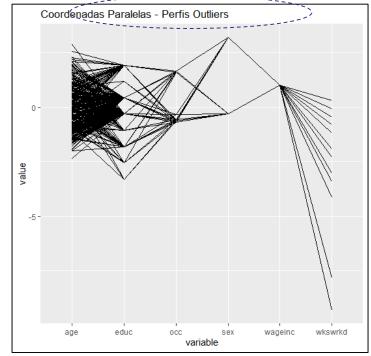


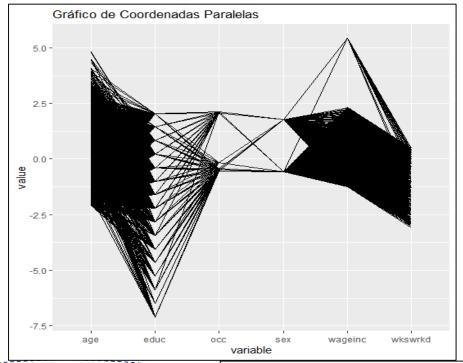
Valores extremos (outliers)

Visualização de partições dos dados:





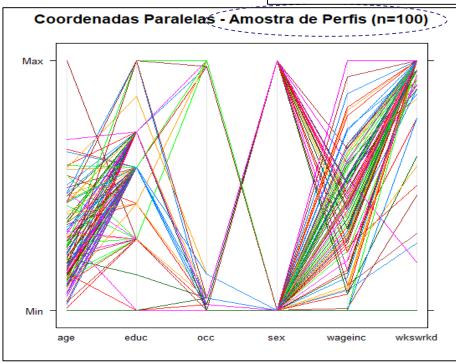




Visualização de partições dos dados:



Duas Amostras aleatórias de 100 perfis



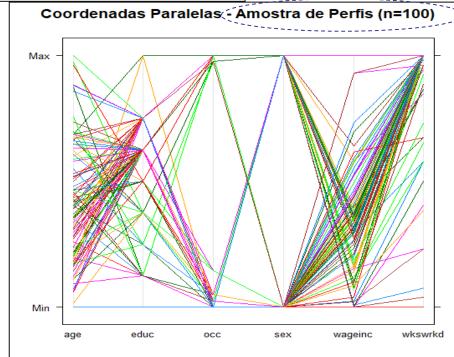
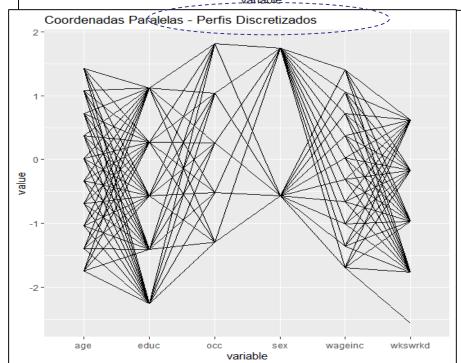
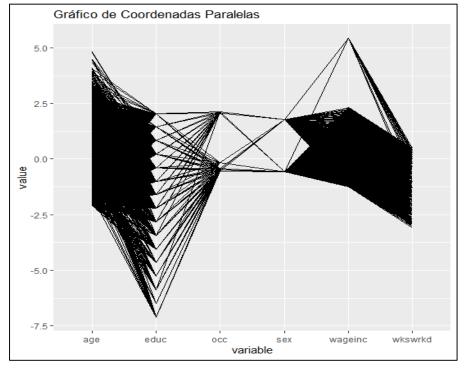


Gráfico de Coordenadas Paralelas 5.0 2.5 -0.0 value -2.5 --5.0 --7.5 age educ wageinc wkswrkd

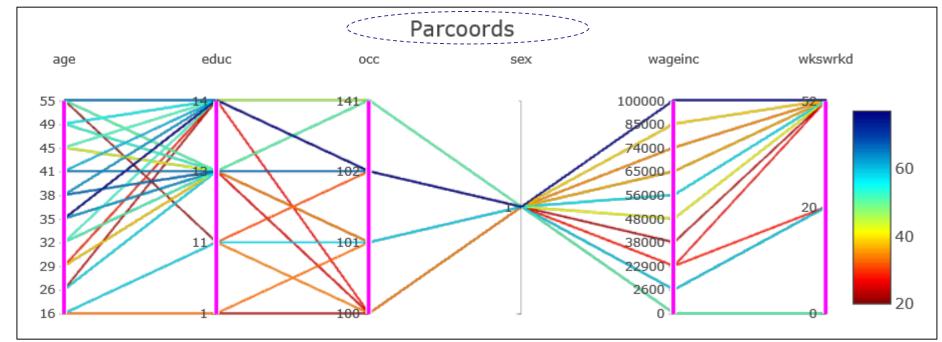
Visualização dos Dados Discretizados:



Discretizar:
Dados
quantitativos
foram
categorizados



Top 100 Perfis mais Frequentes



Dados Multivariados – Big Data

Principais Classes de Problemas (Galeano e Peña, 2019):

- Quantificar Dados provenientes de diferentes fontes (textos, imagens, etc.) e Realizar a Anotação
- Integrar Dados heterogêneos (diferentes escalas)
- Visualização de Dados em Alta Dimensão
- Teste de Múltiplas Hipóteses (α corrigido)
- Métodos automatizados de Seleção de Modelos
- Ajustes de Modelos via soluções Esparsas
- Teoria de Redes (Grafos para Resumir o Padrão de Relação entre Variáveis)

MAE5776-Análise Multivariada

Objetivo: Fornecer base teórica de inferência estatística e fatoração de matrizes na análise de dados multivariados. Apresentar e discutir técnicas de redução de dimensionalidade, integração de bancos de dados e **aprendizado de estruturas**, sob contextos supervisionados e não-supervisionados. Serão abordados casos clássicos (n>p; observações independentes) e soluções em espaços estruturados (n<<p; n muito grande; **observações dependentes**)

Conteúdo (geral):

- 1. Introdução: estrutura de dados, medidas resumo multivariadas, propriedades em espaços duais.
- 2. Distribuição Normal Multivariada: propriedades, estimação, distribuições amostrais, testes de hipóteses para vetores de médias e matrizes de covariância.
 - Elipses de concentração, regiões de confiança, outliers multivariados, gráficos multivariados.
- 3. Técnicas (clássicas) de redução da dimensionalidade (n>p e observações independentes). Teoria de Fatoração de Matrizes na redução de dimensionalidade e integração de bancos de dados.
- 4. Técnicas de redução da dimensionalidade em espaços mais gerais (n<<p>, n muito grande, observações não independentes): soluções em espaços duais, soluções regularizadas e penalizadas, reamostragem.
- 5. **Temas adicionais**: Modelos de Equações Estruturais; Teoria de Grafos Probabilísticos e Aprendizado de Estruturas; Análise de Dados heterogêneos.

Componentes Principais em Dados Heterogêneos

$$Y_{n \times p} = U_{n \times n} \Lambda_r^{1/2} V_{p \times r} \Longrightarrow \min_{a,b} \sum_{i=1}^n \sum_{j=1}^p (Y_{ij} - a_i b_j)^2$$

A: matriz de Escores

B: Matriz de Cargas

GRLM (Generalized Low Rank Model): Extensão da Análise de Componentes Principais (Udell et al., 2016)

$$\min_{a_i,b_j} \sum_{(i,j)\in\Omega} L\Big(Y_{ij} - a_ib_j\Big) + \lambda \sum_{i=1}^n r\Big(a_i\Big) + \lambda \sum_{j=1}^p \tilde{r}\Big(b_j\Big) \quad \text{Obter Arquétipos: "Escores = a_i" e "Cargas = b_j"}$$

Funções Perda Regularizadores

Flexibilidade da solução: combinar diferentes funções Perda e Regularizadores ((L1, L2), bem como realizar a imputação de dados

Modelo	$L_j(u,v)$	r(a)	$\tilde{r}(b)$
CP	$(u-v)^2$	0	0
CPs Regularizados	$(u-v)^2$	$ a _2^2$	$ b _2^2$
NNMF	$(u-v)^2$	$\mathbb{I}(a \ge 0)$	$\mathbb{I}(b \ge 0)$
CP esparso	$(u-v)^2$	$ a _{1}$	$ b _1$
CP Robusto	u-v	$ a _2^2$	$ b _2^2$
CP Logístico	log(1 + esp(-vu))	$ a _2^2$	$ b _2^2$
CP Binário	$(1-vu)_+$	$ a _2^2$	$ b _2^2$
K-means	$(u-v)^2$	$\mathbb{I}(card(a)=1)$	0

Imputação em Dados Missing:

$$\hat{Y}_{ij} = a_i b_j; \quad (i, j) \notin \Omega$$

$$\hat{Y}_{ij} = \arg\min_{x} L_{ij} (a_i b_j, y); \quad (i, j) \notin \Omega$$

Implementação no Python (glrm.py), Julia and H2O.

Ana G Vasconcelos (Mestrado-IME, 2020)