



Universidade de Brasília

DEPARTAMENTO DE ESTATÍSTICA

11 abril 2023

Lista 1

Prof. Dr. George von Borries

Análise Multivariada 1

Aluno: Bruno Gondim Toledo | Matrícula: 15/0167636

1) Exercício de Leitura:

Ler os seguintes artigos:

- Lozano, T. (2016) Data Science: an Emerging Discipline, CITAS.
- Kiliç, G.G.; Turanli, M.; Özden, Ü.H. (2019) Data Mining and Statistics in Data Science, International Social Science Studies Journal.

Os artigos acima são encontrados em <https://www.semanticscholar.org/paper/DATA-MINING-AND-STATISTICS%3A-WHAT%27S-THE-CONNECTIONFriedman/6b3cb8325ca437b6b005027e2e2c64adb99609c8>

- Zola, A. (2019) Machine Learning Engineer vs. Data Scientist. Disponível em <https://www.springboard.com/blog/data-science/machine-learning-engineer-vs-data-scientist/>

Fazer um resumo de três parágrafos e não mais de uma página indicando semelhanças e diferenças entre Estatística, Mineração de Dados e Ciência de Dados. Você pode incluir referência a textos mais recentes sobre estes temas.

O Artigo Lozano (2016) cita que alguns estatísticos vêem a estatística mais como uma ferramenta multidisciplinar e que deve se ater à cálculos probabilísticos para inferência do que propriamente para ser resignificada como ciência de dados. Citando Friedman (1998), diz que mineração de dados (*data mining*) seria mais especificamente sobre encontrar padrões e correspondências nos dados, com foco em grandes bases de dados. Com citação à Dhar (2013), diz que a grande diferença entre estatística e ciência de dados é o volume da informação. Enquanto estatísticos estão acostumados com um volume de dados que caberia em uma calculadora de bolso, a ciência de dados trabalharia com volume de dado gigantesco, necessitando de grande poder computacional. Cita a limitação de algumas técnicas estatísticas, que tendem a serem sensíveis demais com um volume de dados muito grande (*upscaling problem*). Mas mais importante, cita a comum baixa capacidade de estatísticos de lidarem com dados não numéricos e/ou não estruturados, que compoem grande parte da massa de dados produzido ultimamente. Para este problema, a engenharia de dados é importante para estruturar estes dados e torná-los trabalháveis. O texto do Kiliç (2019) foca mais em detalhar e descrever as técnicas estatísticas utilizadas pela ciência de dados para construção de modelos (PCA, regressão, clusterização e afins), bem como a discussão sobre a aplicabilidade dessas técnicas em grandes bases de dados. Ao fim do artigo, os autores tentam traçar correlações entre estatística e ciência de dados. O último artigo Zola (2019), traz mais uma abordagem mercadológica dos termos e as necessidades do mercado quanto a qualificação dos profissionais que pretendem trabalhar em áreas relacionadas a dados.

Ciencia de dados é grande volume de dados. Estatística tem as técnicas mais sofisticadas, mas por vezes anda devagar em relação as transformações da realidade do exponencial crescimento do poder computacional e das massas de dados. A engenharia de dados é a ciência que transforma um dado desestruturado em algo “trabalhável” com técnicas estatísticas numéricas, enquanto a mineração de dados foca em encontrar dentro de uma grande massa de dados, os fragmentos mais relevantes para uma análise.

A grande semelhança é que no fundo essas ciências são quase a mesma coisa e o profissional que exerce alguma dessas funções têm capacidade de compreender e até realizar as outras ciências correlatas; inclusive, num ambiente pequeno-corporativo, a realidade é que provavelmente será o mesmo profissional à executar todos esses serviços. Em Kiliç (2019), é citado que “O ponto em comum da estatística com a ciência de dados é aprender com dados, ou transformar dados em conhecimento” (Kuonen, 2004 apud Kiliç, 2019)

2) Escolha uma área de pesquisa de interesse. Pesquise artigos publicados em revista indexadas e descreva (resumidamente) um exemplo incando o tipo de problema (ou problemas) entre os listados abaixo. Inclua referência bibliográfica e indique as características dos dados e estudo que relacionam ao tipo de problema (ou problemas) indicado.

a) Análise multivariada clássica ($n < p$).

Análises fatoriais e escalas psicométricas. “Com o uso dessa técnica, o investigador discerne os aspectos isolados da estrutura latente em uma matriz de dados e, logo depois, indica o quanto que cada variável é explicada por cada dimensão; ou seja, ele é capaz de descrever as correlações entre muitas variáveis observáveis em termos de um número menor de variáveis não observáveis (linearmente relacionadas com as variáveis originais observadas) denominadas dimensões ou variáveis latentes” (VASCONCELOS, Rafaela Balliana Bragança de. (2016) Apud Hair Jr. et al., 2009; Johnson e Wichern, 2007)[1]

b) Mineração de dados (*Data Mining*) (n elevado).

c) Aprendizado estatístico (*Statistical Learning*) (p elevado).

d) Reconhecimento de padrões.

e) *Data Science*.

Referências:

[1] VASCONCELOS, R. **Avaliação psicométrica de um instrumento para mensuração de autoeficácia**. Monografia (Bacharelado em Estatística) – Departamento de Estatística, Universidade de Brasília. Brasília, p. 21. 2016.

[2]

[3]

[4]

[5]

3) Descreva cada um dos problemas encontrados na análise multivariada, encontre um exemplo de caso real e indique uma falha nas técnicas estatísticas tradicionais.

a) *Mining* (mineração, n muito elevado).

Data mining é a mineração de grandes massas de dados, afim de encontrar correlações e padrões de interesse.

O principal problema ocorrido nesses casos é a não identificação do *software* correto para trabalhar com grandes massas de dados. O *R* por exemplo, por suas especificades de programação pouco eficiente e de ser uma linguagem orientada a objetos em memória *RAM*, não costuma lidar bem com volumes de dados $>5\text{Gb}$. Além disso, testes como o χ^2 tendem a serem sensíveis demais quando $n \rightarrow \infty$, rejeitando h_0 sob pequenas flutuações.

b) *Scalability* (escalabilidade).

A escalabilidade tende a ter um problema parecido com a mineração de dados, visto que pode esbarrar na sensibilidade de alguns testes para grandes massas de dados. Além disso, como a ciência estatística intrinsecamente trabalha com erro (leia-se; variações, variância), por vezes um resultado encontrado não pode ser replicado exatamente, não garantindo assim a escalabilidade automática.

c) *High Dimensional Data* (dados em alta dimensão, $n > p$).

Esse talvez seja um dos tópicos mais difíceis para a estatística, ou seja, trabalhar com conjuntos de poucas observações mas muitas variáveis pode ser um desafio no sentido de não haver parâmetros sólidos a serem buscados, além da dificuldade de encontrar padrões e correlações robustas quando não há muito dado para se comparar.

d) Pequenas amostras.

O principal diferencial da ciência estatística é se aprimorar em dominar a inferência a partir de amostras. Muitas de nossas técnicas são extremamente sofisticadas e teoricamente carregadas justamente a fim de ser possível analisar evidências populacionais sobre pequenas amostras. Porém, quando $n \rightarrow 0$, diversos testes estatísticos como os de normalidade (*Kolmogorov-Smirnov*, *Lilliefors*, etc) tendem a ser muito conservadores com poucos dados; ou seja, não rejeitam h_0 em situações duvidosas.

e) *Curse of Dimensionality* (Problema de dimensionalidade).

Similar aos dados em alta dimensão, o problema da dimensionalidade costuma ocorrer quando são estudadas diversas variáveis em uma amostra; o que compromete a replicabilidade do estudo em uma outra amostra de igual tamanho visto que a quantidade de características latentes por vezes torna inviável a homogeneização dessas amostras.

Em análise fatorial, por exemplo, é muito importante a aplicação de técnicas como a análise de componentes principais, entre outras técnicas, com interesse de reduzir a dimensionalidade do conjunto de dados.

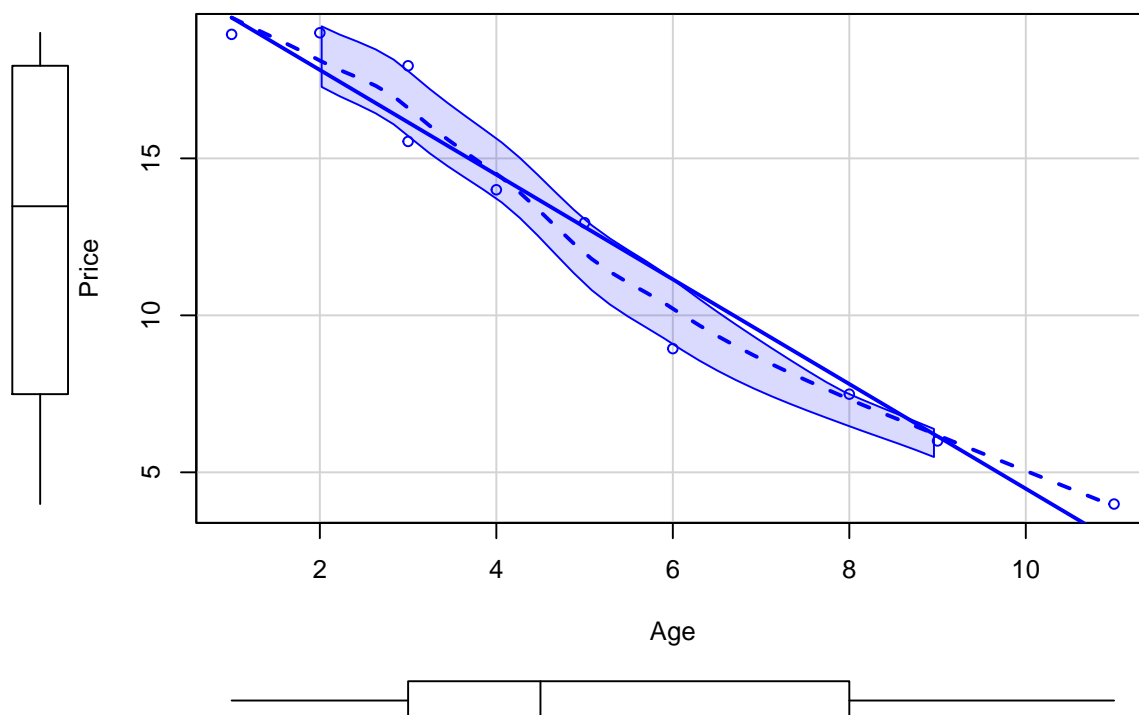
Conforme citado em aula pelo Prof. von Borries, por vezes conjuntos k dimensionais com $k \rightarrow \infty$ contém uma quantidade ϵ de dimensões que correspondem a 80, 90% da informação, sendo o restante responsável pela menor parte. É trabalho do estatístico analisar e tratar estes conjuntos a fim de tornar a análise trabalhável e replicável dado características latentes de real interesse.

4) Fazer os seguintes exercícios do capítulo 1 de *Johnson e Wichern*: 1.2, 1.6, 1.14 e 1.22.

1.2) A mourning newspaper list the following used-car prices for a foreign compact with age (Age) measured in years and selling price (Price) measured in thousand of dollars:

Age	1.00	2	3.00	3.00	4	5.00	6.00	8.00	9	11.00
Price	18.95	19	17.95	15.54	14	12.95	8.94	7.49	6	3.99

a) Construct a scatter plot of the data and marginal dot diagrams.



b) Infer the sign of the sample covariance s_{12} from the scatter plot.

Baseado no gráfico de dispersão, a covariância desses dados é alta e negativa.

A covariância para uma amostra é definida como:

$$s_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

A covariância calculada para este conjunto de dados é:

[1] -17.71

O que confirma o palpite inferido sobre o gráfico de dispersão.

- c) Compute the sample means \bar{x}_1 and \bar{x}_2 and the sample variances s_{11} and s_{22} . Compute the sample covariance s_{12} and the sample correlation coefficient r_{12} . Interpret these quantities.

	Age	Price
Média	5.20	12.48
Variância	10.62	30.85

A covariância da amostra s_{12} observada foi -17.71; enquanto o coeficiente de correlação da amostra r_{12} observado foi -0.98. Esta covariância é visualmente analisável pelo gráfico de dispersão conforme questão anterior, enquanto o valor observado na correlação indica que ela é muito forte à esquerda, quase absoluta.

- d) Display the sample mean array \bar{x} , the sample variance-covariance array S_n , and the sample correlation array \mathbf{R} using (1-8).

1.6)

1.14)

1.22)