

EdX and its Members use cookies and other tracking technologies for performance, analytics, and marketing purposes. By using this website, you accept this use. Learn more about these technologies in the [Privacy Policy](#).



Course **Discussion** Progress Resources

☰ All Topics

Add a Post

Search all posts

Search

Filter Topics

filter topics ▼

All Discussions

★ Posts I'm Following

(Optional) Unit 8 Principal component analysis:(Optional) Lecture 23: Principal Component Analysis

1. Objectives

2. Introduction

3. Multivariate Statistics and Geometry Behind the Empirical Covariance

4. Geometry Behind the Empirical Covariance Matrix - Continued

5. Review of Linear Algebra Required for this Lecture

Resolution of the Final Exam in Julia

discussion posted 2 days ago by [Lobianco](#)

The resolution of Final Exam (except theoretical-only questions) in Julia (copyable code at the bottom) :



6. Principal Component Analysis
(PCA) - Theorem

7. Largest Eigenvalue and Principal
Directions



```

1 # Final Exam of MITx 18.6051x, December 2019
2
3 using Distributions, SymPy, LinearAlgebra, CSV, DataFrames, TableView
4
5 # *****
6 # ** Page 2 - Review of fundamentals *****
7 # *****
8
9 # *****
10 # * Square of a standard normal: Warmup
11 # Just the Chi Square with 1 degree of freedom
12
13 # *****
14 # * Approximation via Central Limit Theorem
15
16  $\chi_1 = \text{Chisq}(1)$  # T-distribution  $\{\text{Chisq}\{\text{Float64}\}(v=1.0)\}$ 
17  $\alpha = 0.05$   $\{0.0500\}$ 
18  $q = \text{quantile}(\chi_1, 1-\alpha)$   $\{3.84\ldots\}$ 
19
20 # Method #1: using CLT:
21  $A_{100} = \text{Normal}(1, \sqrt{2/100})$   $\{\text{Normal}\{\text{Float64}\}(\mu=1.0, \sigma=0.1414213562373095)\}$ 
22  $A_{16} = \text{Normal}(1, \sqrt{2/10^6})$   $\{\text{Normal}\{\text{Float64}\}(\mu=1.0, \sigma=0.001414213562373095)\}$ 
23  $q = \text{quantile}(A_{100}, 1-\alpha)$   $\{1.23\ldots\}$ 
24  $q = \text{quantile}(A_{16}, 1-\alpha)$   $\{1.00\ldots\}$ 
25
26 # Method 2: true values using xhi-squared distribution without the CLT
27 # approximation:
28  $\chi_{100} = \text{Chisq}(100)$   $\{\text{Chisq}\{\text{Float64}\}(v=100.0)\}$ 
29  $q = \text{quantile}(\chi_{100}, 1-\alpha)/100$   $\{1.24\ldots\}$ 
30  $\chi_{16} = \text{Chisq}(10^6)$   $\{\text{Chisq}\{\text{Float64}\}(v=1.0e6)\}$ 
31  $q = \text{quantile}(\chi_{16}, 1-\alpha)/10^6$   $\{1.00\ldots\}$ 
32
33 # *****
34 # ** Page 3 - Exponential family *****
35 # *****
36
37 # *****
38 # * Canonical form
39  $\theta = \text{symbols}(\theta, \text{real}=\text{true}, \text{negative}=\text{true})$   $\{\theta\}$ 
40  $\mu, x, \Sigma, \bar{x}, n = \text{symbols}(\mu \ x \ \Sigma \ \bar{x} \ n, \text{positive}=\text{true}, \text{real}=\text{true})$   $\{(\mu, x, \Sigma, \bar{x}, n)\}$ 
41 # Original function..
42  $f\mu = (1/\text{sympy.sqrt}(2*\text{sympy.pi}*x^3))*\text{sympy.exp}(-(x-\mu)^2/(2*\mu^2*x))$   $\{\text{sqrt}(2)*\text{exp}(-(x-\mu)^2/(2*x*\mu^2))/(2*\text{sqrt}(\text{pi})*x^{(3/2)})\}$ 
43 # Canonical exponential form..
44  $\theta\text{Expr} = -1/(2*\mu^2)$   $\{-1/(2*\mu^2)\}$ 
45  $b\theta\text{Expr} = -\text{sympy.sqrt}(-2*\theta)$   $\{-\text{sqrt}(2)*\text{sqrt}(-\theta)\}$ 
46  $c\theta\text{Expr} = -1/(2*x) + \text{sympy.log}(1/\text{sympy.sqrt}(2*\text{sympy.pi}*x^3))$   $\{\log(\text{sqrt}(2)/(2*\text{sqrt}(\text{pi})*x^{(3/2)})) - 1/(2*x)\}$ 
47  $f\theta = \text{exp}(x*\theta - b\theta\text{Expr} + c\theta\text{Expr})$   $\{\text{sqrt}(2)*\text{exp}(x*\theta + \text{sqrt}(2)*\text{sqrt}(-\theta) - 1/(2*x))/(2*\text{sqrt}(\text{pi})*x^{(3/2)})\}$ 
48
49 # Checks the two versions are the same

```

```

50 fμ2 = subs(fθ,Dict(θ => θExpr)) [sqrt(2)*exp(-x/(2*μ^2) + 1/μ - 1/(2*x))/(2*sqrt(pi)*x^(3/2))
51 test = simplify(fμ - fμ2) [0]
52
53 # *****
54 # * Canonical link
55 gExpr = solve(Eq(θExpr,θ),μ)[1] [sqrt(2)/(2*sqrt(-θ))]
56
57 # *****
58 # * Expectation and variance
59 Eθ = diff(bθExpr,θ) [-sqrt(2)*sqrt(-θ)/(2*θ)]
60 Eμ = subs(Eθ,Dict(θ => θExpr)) [μ]
61 Varθ = diff(Eθ,θ) [sqrt(2)*sqrt(-θ)/(4*θ^2)]
62 Varμ = subs(Varθ,Dict(θ => θExpr)) [μ^3]
63
64 # *****
65 # * Fisher Information
66 l1 = sympy.log(fθ) [log(sqrt(2)*exp(x*θ + sqrt(2)*sqrt(-θ) - 1/(2*x))/(2*sqrt(pi)*x^(3/2)))]
67 dl1 = simplify(diff(l1,θ)) [x + sqrt(2)*sqrt(-θ)/(2*θ)]
68 dl2 = simplify(diff(dl1,θ)).as_real_imag()[1] [sqrt(2)/(4*(-θ)^(3/2))]
69 FisherIθ = simplify(subs(-dl2,Dict(x=>Eθ))) [-sqrt(2)/(4*(-θ)^(3/2))]
70 FisherIμ = subs(FisherIθ,Dict(θ => θExpr)) # Still I of θ, even if expressed in terms of μ [-μ^3]
71
72 # *****
73 # * MLE
74 Lnθ = exp(θ*Σx-n*bθExpr) # omitting c(x,phi) [exp(sqrt(2)*n*sqrt(-θ) + Σx*θ)]
75 lnθ = log(Lnθ) [sqrt(2)*n*sqrt(-θ) + Σx*θ]
76 dlnθ = diff(lnθ,θ) [sqrt(2)*n*sqrt(-θ)/(2*θ) + Σx]
77 θhat = subs(solve(dlnθ, θ)[1], n/Σx => 1/χ) [-1/(2*χ^2)]
78
79 Lnμ = exp(-Σx/(2*μ^2)+n/μ) # omitting c(x,phi) [exp(n/μ - Σx/(2*μ^2))]
80 lnμ = log(Lnμ) [n/μ - Σx/(2*μ^2)]
81 dlnμ = diff(lnμ,μ) [-n/μ^2 + Σx/μ^3]
82 μhat = subs(solve(dlnμ, μ)[1], Σx/n => x̄) [x̄]
83
84 # *****
85 # * Asymptotic variance
86 # We already have the variance of θhat, but we don't have those of μhat
87 L1μ = exp(-x/(2*μ^2)+1/μ) # omitting c(x,phi) [exp(-x/(2*μ^2) + 1/μ)]
88 l1μ = log(L1μ) [-x/(2*μ^2) + 1/μ]
89 dl1μ = diff(l1μ,μ) [x/μ^3 - 1/μ^2]
90 dl2μ = diff(dl1μ,μ) [-3*x/μ^4 + 2/μ^3]
91 FisherIμ2 = simplify(subs(-dl2μ,Dict(x=>Eμ))) # This is the Fisher information of μ, expressed in terms of μ [μ^(-3)]
92 Varμ = 1/FisherIμ2 [μ^3]

```

```

94 # *****
95 # ** Page 4 - Voting regression *****
96 # *****
97
98 # *****
99 # * Covariance matrix
100 # Just checking which option result in a scalar
101 A=rand(5,5) |> 5x5 Array{Float64,2}:
102 u = [0,0,0,-1,-1] |> Vector{Int64} with 5 elements
103 d4 = u'*A*u |> 2.24...
104 d5 = u*A*u' |> DimensionMismatch("matrix A has dimensions (5,1), matrix B has dimensions (5,5)")
105
106 # *****
107 # * Estimate the Variance
108 cd(@_DIR_) | ✓
109 df = CSV.read("data_gerber_trunc.csv",delim=",") |> 20000x16 DataFrame. Omitted printing of 8 columns
110 showtable(df) |> WebIO.Scope
111 n = size(df)[1] |> 20000
112 X = hcat(ones(n),df.civicduty,df.hawthorne,df.self,df.neighbors) |> 20000x5 Array{Float64,2}:
113 Y = df.voting |> 20000-element LazyArrays.ApplyArray{Int64,1,typeof(vcat),NTuple{4,CSV.Column{Int64,Int64}}}:
114 β = (X' * X)^-1 * X' * Y |> Vector{Float64} with 5 elements
115 ε = Y-X*β # Residuals |> 20000-element LazyArrays.ApplyArray{Float64,1,typeof(vcat),NTuple{4,Array{Float64,1}}}:
116 σ² = sum([εᵢ² for εᵢ in ε])/(n-1) |> 0.216...
117 Σ = σ² * (X' * X)^-1 |> 5x5 Array{Float64,2}:
118 u = [0,0,0,-1,-1] |> Vector{Int64} with 5 elements
119 Varβ4minusβ3 = u' * Σ * u |> 0.000272...

```

```

121 # *****
122 # ** Page 5 - Independence tests for Bernoulli / regression *****
123 # *****
124
125 r, p, q, = symbols("r p q", positive=True, real=True) (r, p, q)
126
127 # The key in this exercise was to recognise that  $XY \sim \text{Ber}(r)$ , that is,  $XY = 1$  if and only if  $X$  and  $Y$  are both 1
128
129 # *****
130 # * Asymptotic Variance Under the Null
131
132  $\Sigma = \begin{bmatrix} p(1-p) & r-pq & r(1-p) \\ r-pq & q(1-q) & r(1-q) \\ r(1-p) & r(1-q) & r(1-r) \end{bmatrix}$  > 3x3 Array{Sym,2}:
133  $\omega = [-q, -p, 1]$  > Vector{Sym} with 3 elements
134
135 # note that even under the null the estimators haven't zero correlation, i.e using
136 #  $\Sigma_\theta = \begin{bmatrix} p(1-p) & 0 & 0 \\ 0 & q(1-q) & 0 \\ 0 & 0 & p q (1-p q) \end{bmatrix}$  would be wrong
137  $V_\theta = \text{simplify}(\omega' \cdot \Sigma \cdot \omega)$   $-4p^2q^2 + p^2q + pq^2 + 6pq^2r - 2p^2r - 2q^2r - r^2 + r$ 
138  $\#V_\theta = \text{simplify}(\text{subs}(V_\theta, \text{Dict}(r \Rightarrow pq)))$   $pq(-3pq + p + q + 1)$ 
139  $V_\theta = \text{simplify}(\text{subs}(V_\theta, \text{Dict}(r \Rightarrow pq)))$   $pq(pq - p - q + 1)$ 
140
141 # *****
142 # * Happiness and Being in a Relationship
143  $n = 205 + 179 + 301 + 315$  1000
144  $\text{hatp} = (205 + 179) / (205 + 179 + 301 + 315)$  0.384
145  $\text{hatq} = (205 + 301) / (205 + 179 + 301 + 315)$  0.506
146  $\text{hatr} = (205) / (205 + 179 + 301 + 315)$  0.205
147
148  $v0num = N(V_\theta.\text{evalf}(\text{subs}=\text{Dict}(p \Rightarrow \text{hatp}, q \Rightarrow \text{hatq}, r \Rightarrow \text{hatr})))$  0.0591...
149  $Tn = \text{abs}(\text{sqrt}(n) * (\text{hatr} - \text{hatq} * \text{hatp}) / \text{sqrt}(v0num))$  1.39...
150  $\alpha = 0.05$  0.0500
151  $q = \text{quantile}(\text{Normal}(), 1 - (\alpha/2))$  1.96...
152  $Tn > q$  false
153  $p\_value = 2 * \text{cdf}(\text{Normal}(), -Tn)$  0.1642258510754747
154

```

Copyable code:

```

# Final Exam of MITx 18.6051x, December 2019

using Distributions, SymPy, LinearAlgebra, CSV, DataFrames, TableView

# *****
# ** Page 2 - Review of fundamentals *****
# *****

# *****
# * Square of a standard normal: Warmup
# Just the Chi Square with 1 degree of freedom

# *****
# * Approximation via Central Limit Theorem

 $\chi_1$  = Chisq(1) # T-distribution
 $\alpha$  = 0.05
q = quantile( $\chi_1$ , 1- $\alpha$ )

# Method #1: using CLT:
A100 = Normal(1, sqrt(2/100))
A16 = Normal(1, sqrt(2/10^6))
q = quantile(A100, 1- $\alpha$ )
q = quantile(A16, 1- $\alpha$ )

# Method 2: true values using xhi-squared distribution without the CLT
# approximation:
 $\chi_{100}$  = Chisq(100)
q = quantile( $\chi_{100}$ , 1- $\alpha$ )/100
 $\chi_{16}$  = Chisq(10^6)
q = quantile( $\chi_{16}$ , 1- $\alpha$ )/10^6

# *****
# ** Page 3 - Exponential family *****
# *****

# *****
# * Canonical form
 $\theta$  = symbols(" $\theta$ ", real=true, negative=true)
 $\mu$ , x,  $\Sigma$ x,  $\bar{x}$ , n = symbols(" $\mu$  x  $\Sigma$   $\bar{x}$  n", positive= true, real=true)
# Original function..
f $\mu$  = (1/sympy.sqrt(2*sympy.pi*x^3))*sympy.exp(-(x- $\mu$ )^2/(2* $\mu$ ^2*x))
# Canonical exponential form..
 $\theta$ Expr = -1/(2* $\mu$ ^2)

```

```

b0Expr = -sympy.sqrt(-2*θ)
c0Expr = -1/(2*x)+sympy.log(1/sympy.sqrt(2*sympy.pi*x^3))
f0 = exp(x*θ-b0Expr+c0Expr)

# Checks the two versions are the same..
fμ2 = subs(f0,Dict(θ => θExpr))
test = simplify(fμ - fμ2)

# *****
# * Canonical link
gExpr = solve(Eq(θExpr,θ),μ)[1]

# *****
# * Expectation and variance
Eθ = diff(b0Expr,θ)
Eμ = subs(Eθ,Dict(θ => θExpr))
Varθ = diff(Eθ,θ)
Varμ = subs(Varθ,Dict(θ => θExpr))

# *****
# * Fisher Information
l1 = sympy.log(f0)
dl1 = simplify(diff(l1,θ))
dl2 = simplify(diff(dl1,θ)).as_real_imag()[1]
FisherIθ = simplify(subs(-dl2,Dict(x=>Eθ)))
FisherIμ = subs(FisherIθ,Dict(θ => θExpr)) # Still I of θ, even if expressed in terms of μ

# *****
# * MLE
Lnθ = exp(θ*Σx-n*b0Expr) # omitting c(x,phi)
lnθ = log(Lnθ)
dlnθ = diff(lnθ,θ)
θhat = subs(solve(dlnθ, θ)[1], n/Σx => 1/x)

Lnμ = exp(-Σx/(2*μ^2)+n/μ) # omitting c(x,phi)
lnμ = log(Lnμ)
dlnμ = diff(lnμ,μ)
μhat = subs(solve(dlnμ, μ)[1], Σx/n => x̄)

# *****
# * Asymtotic variance
# We already have the variance of θhat, but we don't have those of μhat
L1μ = exp(-x/(2*μ^2)+1/μ) # omitting c(x,phi)
l1μ = log(L1μ)

```



```

dl1μ = diff(l1μ,μ)
dl2μ = diff(dl1μ,μ)
FisherIμ2 = simplify(subs(-dl2μ,Dict(x=>Eμ))) # This is the Fisher information of μ, expressed in terms of μ
Varμ = 1/FisherIμ2

# *****
# ** Page 4 - Voting regression *****
# *****

# *****
# * Covariance matrix
# Just checking which option result in a scalar
A=rand(5,5)
u = [0,0,0,-1,-1]
d4 = u'*A*u
d5 = u*A*u'

# *****
# * Estimate the Variance
cd(@__DIR__)
df = CSV.read("data_gerber_trunc.csv",delim=",")
showtable(df)
n = size(df)[1]
X = hcat(ones(n),df.civicduty,df.hawthorne,df.self,df.neighbors)
Y = df.voting
β = (X' * X)^-1 * X' * Y
ε = Y-X*β # Residuals
σ² = sum([ε_i^2 for ε_i in ε])/(n-1)
Σ = σ² * (X' * X) ^-1
u = [0,0,0,-1,-1]
Varβ4minusβ3 = u' * Σ * u

# *****
# ** Page 5 - Independence tests for Bernoulli / regression *****
# *****

r, p, q,= symbols("r p q", positive=true, real=true)

# The key in this exercise was to recognise that  $XY \sim \text{Ber}(r)$ , that is,  $XY = 1$  if and only if  $X$  and  $Y$  are both 1

# *****
# * Asymptotic Variance Under the Null

```

```

Σ = [ p*(1-p)  r-p*q  r*(1-p) ; r-p*q  q*(1-q)  r*(1-q) ; r*(1-p)  r*(1-q)  r*(1-r)]
ω = [-q, -p, 1]

# note that even under the null the estimators haven't zero correlation, i.e using
# Σ₀ = [ p*(1-p)  0  0 ; 0  q*(1-q)  0 ; 0  0  p*q*(1-p*q)] would be wrong
V₀ = simplify(ω' * Σ * ω)
#V₀ = simplify(subs(V₀ , Dict(r => p*q)))
V₀ = simplify(subs(V₀ , Dict(r => p*q)))

# *****
# * Happiness and Being in a Relationship
n = 205+179+301+315
hatp = (205+179)/(205+179+301+315)
hatq = (205+301)/(205+179+301+315)
hatr = (205)/(205+179+301+315)

v0num = N(V₀.evalf(subs=Dict(p => hatp, q => hatq, r => hatr)))
Tn = abs(sqrt(n)*(hatr-hatq*hatp)/sqrt(v0num))
α = 0.05
q = quantile(Normal(),1-(α/2))
Tn > q
p_value = 2 * cdf(Normal(),-Tn)

```

Related to: [Final exam:Final Exam / 5.](#)

This post is visible to everyone.

0 responses

Preview

Submit



[About](#)
[edX for Business](#)

[Terms of Service & Honor Code](#)
[Privacy Policy](#)
[Accessibility Policy](#)

[Blog](#)
[Contact Us](#)
[Help Center](#)



© 2019 edX Inc. All rights reserved.
| 深圳市恒宇博科技有限公司 [粤ICP备17044299号-2](#)

