

NOTEBOOK DESTINADO A REGISTRAR AS TAREFAS DA DISCIPLINA DE AEDI - 1º/2025

PROFESSOR: JOÃO GABRIEL DE MORAES SOUZA

ALUNO: GUSTAVO PARREIRA LIMA CUNHA

TAREFA 1

✓ QUESTÃO A

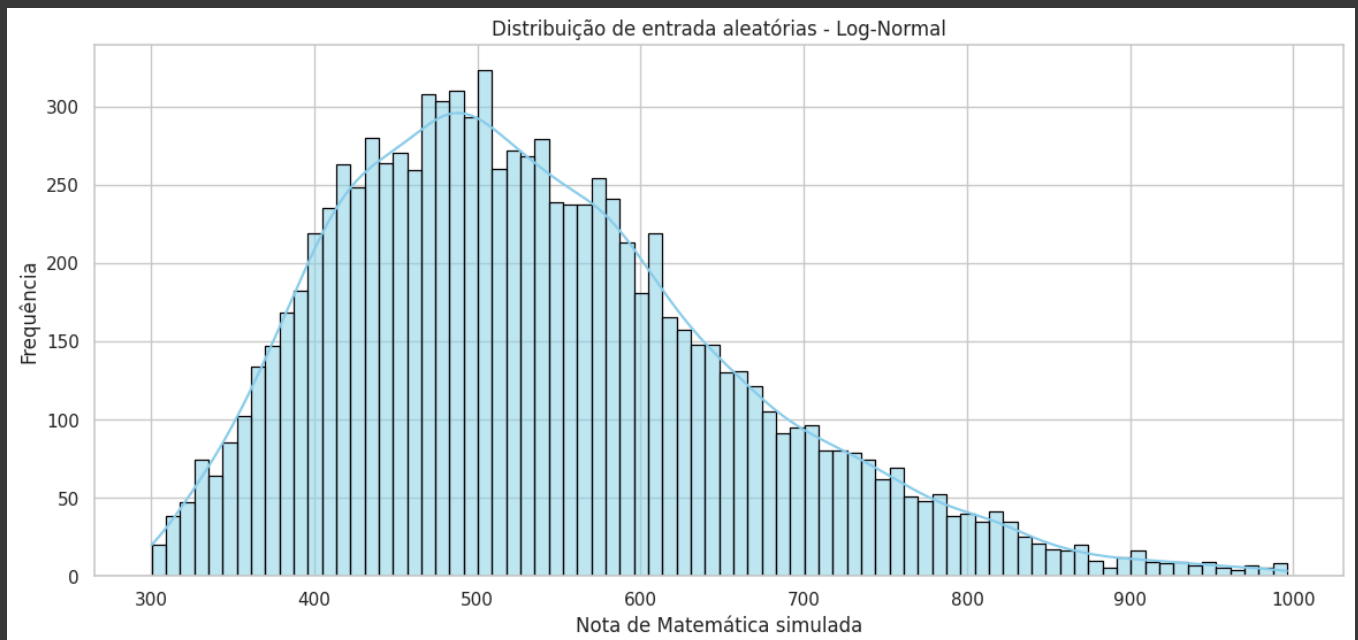
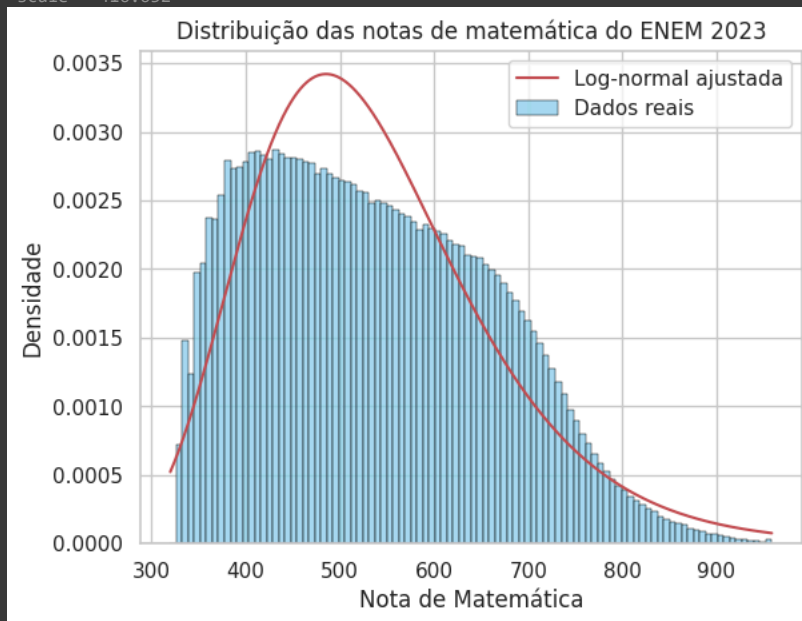
```
1 # Fenômeno representado: Notas obtidas por candidatos do ENEM na prova de
  Matemática
2
3 # Distribuição escolhida: Log-Normal
4
5 # Fonte: Microdados do ENEM 2023
6 # Link: https://download.inep.gov.br/microdados/microdados\_enem\_2023.zip
7
8 # Leitura dos dados CSV
9
10 import pandas as pd
11 import numpy as np
12 import matplotlib.pyplot as plt
13 import seaborn as sns
14 from scipy.stats import lognorm
15
16 # Especificação do caminho do arquivo csv
17
18 path = '/content/drive/MyDrive/AEDI/MICRODADOS_ENEM_2023.csv'
19
20 # Lendo somente a coluna de interesse do arquivo de microdados do ENEM 2023,
  referente à nota da prova de matemática(NU_NOTA_MT), para otimização de
  performance:
21
22 enem = pd.read_csv(path, usecols=['NU_NOTA_MT'], encoding='latin1', sep=';')
23
24 # Remover valores nulos e negativos/zero para cálculo de média e desvio-padrão
25
26 dados = enem['NU_NOTA_MT']
27 dados = dados.dropna()
28 dados = dados[dados > 0]
29
30 # Ajuste dos dados a uma distribuição log-normal
31 shape, loc, scale = lognorm.fit(dados)
32
33 # Calculando  $\mu$  e  $\sigma$  (da normal associada)
34 sigma = shape
35 mu = np.log(scale)
36
37 # Média real da log-normal:
38 u = np.exp(mu + (sigma**2) / 2)
39
40 # Desvio padrão real da log-normal:
41 sd = np.sqrt((np.exp(sigma**2) - 1) * np.exp(2 * mu + sigma**2))
42
43 print(f"Média real da log-normal: {u:.2f}")
44 print(f"Desvio padrão real da log-normal: {sd:.2f}")
45 print(f"Parâmetros da log-normal: \n shape = {round(shape, 3)}\n loc = {round
  (loc, 3)}\n scale = {round(scale, 3)}")
46
47 # Geração do gráfico com o valor das notas de matemática do ENEM 2023:
  Histograma
48
49 # Histograma dos dados reais
50 sns.histplot(dados, bins=100, stat='density', color='skyblue',
  edgecolor='black', label='Dados reais')
51
52 # Gerar valores do eixo x
53 x = np.linspace(min(dados), max(dados), 10000)
54
55 # ajuste manual de parâmetros, para observar resposta da curva teórica (função
  densidade de probabilidade)
56
57 shape_adj = shape*1
```

```
58
59 loc_adj = loc*1
60
61 scale_adj = scale*1
62
63 pdf = lognorm.pdf(x, s=shape_adj, loc=loc_adj, scale=scale_adj)
64
65 # Plot da curva teórica (FDP)
66 plt.plot(x, pdf, 'r-', label='Log-normal ajustada')
67 plt.xlabel('Nota de Matemática')
68 plt.ylabel('Densidade')
69 plt.title('Distribuição das notas de matemática do ENEM 2023')
70 plt.legend()
71 plt.grid(True)
72 plt.show()
73 print()
74
75 # Definição do número de simulações:
76
77 n_sim = 10000
78
79 # Geração das entradas aleatoriamente respeitando-se a distribuição log-normal
80 com a média e desvio-padrão obtidos a partir dos dados reais do Enem 2023.
81
82 entradas = []
83
84 entradas = lognorm.rvs(s=shape, loc=loc, scale=scale, size=n_sim)
85
86 # Normalizando para o intervalo esperado
87
88 min_nota_real = 300
89 max_nota_real = 1000
90
91 # Filtrar as entradas, mantendo apenas as que estão dentro do intervalo válido
92 entradas_filtradas = entradas[(entradas >= min_nota_real) & (entradas <=
93 max_nota_real)]
94
95 # Geração do gráfico dos valores aleatórios seguindo uma distribuição
96 log-normal, com a média e desvio-padrão obtidos dos dados reais.
97
98 plt.figure(figsize=(14, 6))
99 sns.histplot(entradas_filtradas, bins=80, kde=True, color='skyblue',
100 edgecolor='black')
101
102 plt.title('Distribuição de entrada aleatórias - Log-Normal')
103 plt.xlabel('Nota de Matemática simulada')
104 plt.ylabel('Frequência')
105 plt.grid(True)
106 plt.show()
107
```

```

Média real da log-normal: 434.79
Desvio padrão real da log-normal: 129.77
Parâmetros da log-normal:
shape = 0.292
loc = 102.642
scale = 416.632

```



✓ Definição de uma possível pergunta problema:

Considerando que a distribuição das notas de Matemática no ENEM 2023 segue uma distribuição log-normal com média $\mu = 434.79$ e desvio padrão $\sigma = 129.77$, responda a problemática abaixo.

Problemática:

Para que um estudante seja competitivo na disputa por uma vaga no curso de Engenharia de Computação na Universidade de Brasília (UnB), é fundamental que sua nota de matemática esteja no intervalo entre 700 e 850 pontos.

Qual é a proporção esperada de candidatos que, com base no modelo de simulação, conseguiriam atingir uma nota dentro desse intervalo competitivo?

```

1 # # A resposta está na área sob a curva log-normal situada no intervalo entre 700 e 850 pontos.
2
3 # X(nota de matemática estar entre 700 e 850 pontos)

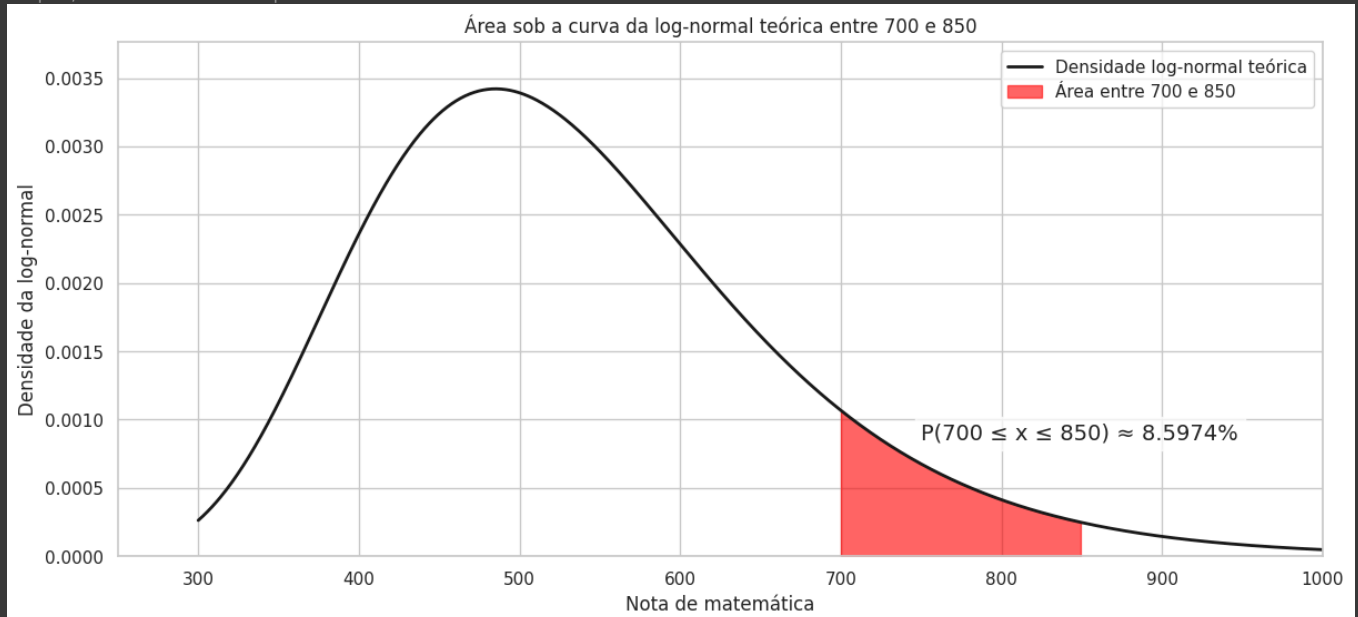
```

```

4 # variável aleatória  $f(x) = X(700 \leq x \leq 850)$ 
5
6 from scipy.stats import lognorm
7
8 # Estimar a área sob a curva teórica:  $P(x \leq 850) - P(x \leq 700)$ 
9 area = lognorm.cdf(850, s=shape, loc=loc, scale=scale) - lognorm.cdf(700, s=shape, loc=loc, scale=scale)
10 print(f"Proporção entre 700 e 850 pela CDF teórica: {100*area:.4f}%")
11
12 # Gerar eixo x e y da curva teórica
13 x = np.linspace(300, 1000, 1000)
14 pdf_teorica = lognorm.pdf(x, s=shape, loc=loc, scale=scale)
15
16 # Máscara para o intervalo do problema
17 mask = (x >= 700) & (x <= 850)
18
19 # Desenho da curva e preenchimento da área
20 plt.figure(figsize=(14, 6))
21 plt.plot(x, pdf_teorica, 'k-', linewidth=2, label='Densidade log-normal teórica')
22 plt.fill_between(x[mask], pdf_teorica[mask], color='red', alpha=0.6, label='Área entre 700 e 850')
23
24 # Ajuste de texto com a resposta teórica
25 plt.text(750, max(pdf_teorica)*0.25, f" $P(700 \leq x \leq 850) \approx \{100*area:.4f\}\%$ ", fontsize=14, bbox=dict(facecolor='white', alpha=0.8))
26
27 # Estética
28 plt.title("Área sob a curva da log-normal teórica entre 700 e 850")
29 plt.xlabel("Nota de matemática")
30 plt.ylabel("Densidade da log-normal")
31 plt.xlim([250, 1000])
32 plt.ylim([0, max(pdf_teorica)*1.1])
33 plt.grid(True)
34 plt.legend()
35 plt.show()
36

```

Proporção entre 700 e 850 pela CDF teórica: 8.5974%



QUESTÃO B

A apresentação detalhada dos fundamentos estatísticos da simulação, com justificativas formais para a escolha da distribuição utilizada, considerando suas propriedades, suposições e aderência ao contexto.

- A simulação tem como objetivo estudar a distribuição das notas dos candidatos na prova de matemática do Enem 2023, o que foi realizado mediante geração de simulações realistas para responder a uma situação problema formulada: "qual a chance de um candidato atingir entre 700 e 850 pontos?"
- Fundamentos da escolha da distribuição: a variável objeto de estudo - nota - é contínua, positiva e possui assimetria. A Teoria de Resposta ao Item (TRI) faz com que não se tenha notas situadas na cauda esquerda da distribuição, tendo início em aproximadamente 300 pontos. Em seguida, nota-se uma subida muito rápida, assim como a log-normal. Há um pico em aproximadamente 400 pontos e em seguida um lento decaimento, caracterizando uma cauda longa à direita, como também se observa na log-normal. A principal diferença observada entre a log-normal teórica e os dados reais obtidos do Enem é no formato do decaimento, sendo mais lento nos dados reais quando comparado à curva teórica. Apesar da diferença, a log-normal foi a distribuição escolhida que melhor representa o perfil de distribuição das notas de matemática do Enem 2023.
- Propriedades da Log-normal: Segundo Droubi et al. (2018), diz-se que uma variável aleatória X tem distribuição log-normal se o seu logaritmo $\ln(X)$ tem distribuição normal.

1. Valor esperado: $E[X] = e^{\mu + (\sigma^2/2)}$

2. Variância: $\text{Var}(X) = e^{(2\mu + \sigma^2)} * (e^{\sigma^2} - 1)$

- Essa distribuição permite representar bem caudas longas, conforme as notas do Enem se apresentam.
- A partir de então, foram obtidos os valores dos parâmetros da curva log-normal por meio da função `lognorm.fit(dados)`, que ajustou a coluna que continha as notas de matemática do Enem 2023 a uma curva log-normal, atribuindo-se valores aos seguintes parâmetros:
 1. Shape: afetado pelo desvio padrão, esse parâmetro determina o formato da curva, e não tem interferência sobre o posicionamento no eixo x ou sobre a altura. No presente caso, `shape = 0.292`
 2. Loc: determina o deslocamento horizontal da curva sobre o eixo x. No presente caso, `loc = 102.642`
 3. Scale: é a mediana dos dados. Relaciona-se à altura da curva. No presente caso, `scale = 416.632`
- A partir dos parâmetros obtidos ao se executar o fit nos dados reais, foram geradas entradas aleatórias utilizando-se os mesmos parâmetros, por meio da função `lognorm.rvs(...)` para calcular a probabilidade específica e responder à questão problema formulada.
- Também foi possível concluir que a distribuição log-normal é uma aproximação, pois não é totalmente aderente ao histograma dos dados reais, que possui cauda longa menos pesada que a curva teórica.

BIBLIOGRAFIA

(1) Droubi, Luiz Fernando Palin, Willian Zonato, and Norberto Hochheim. "Distribuição lognormal: propriedades e aplicações na engenharia de avaliações." Congresso de Cadastro Multifinalitário e Gestão Territorial, Florianópolis, Anais... SC. 2018.

QUESTÃO C

A análise crítica dos resultados obtidos via Simulação de Monte Carlo, com discussão sobre:

1. Como a distribuição escolhida impacta os resultados da simulação;
2. As implicações práticas e possíveis inferências decorrentes;
3. A sensibilidade dos resultados às variações dos parâmetros da distribuição.

--

- A simulação só irá gerar entradas aleatórias legítimas se a distribuição escolhida for adequada para representar o fenômeno. No presente caso, como as principais distorções entre a curva real e teórica se encontram na faixa entre 650 e 750 pontos, com a curva real apresentando mais ocorrências desses valores quando comparado à curva teórica, sabe-se que as entradas aleatórias, que foram geradas a partir dos parâmetros da curva teórica, terão menos ocorrências do que deveria para essa faixa de valores. Quanto maior for a aderência dos dados reais à distribuição teórica, menor será esse erro na geração de entradas aleatórias.
- No presente contexto, foi possível utilizar os dados de toda a população para cálculo dos parâmetros da log normal e geração de