

Documentação Técnica do DataSynth e DataSynthDecrypter

Viviane Viana Sofiste de Abreu

June 21, 2024

Contents

1	Introdução	2
2	Importações e Bibliotecas	2
3	Funções de Geração de Ruído Branco	2
3.1	Transformação de Box-Muller	2
4	Funções de Criptografia	3
4.1	Derivar Chave a partir da Senha	3
4.2	Adicionar Ruído Branco e Encriptar Dados	3
5	Funções de Anonimização	4
5.1	Anonimizar Dados Textuais com Ruído Branco	4
5.2	Gerar Imagem de Ruído Branco	5
6	Funções de Interface do Streamlit	5
6.1	Ajustar Estilo CSS	5
6.2	Formatar Colunas do DataFrame	6
6.3	Gerar Senha Aleatória	6
6.4	Decryptar Dados	6
7	Fluxo Principal do DataSynth	7
7.1	Interface do Streamlit	7
8	Fluxo Principal do DataSynthDecrypter	11
8.1	Interface do Streamlit	11

1 Introdução

O DataSynth é uma ferramenta projetada para anonimizar dados sensíveis de transações financeiras utilizando técnicas de ruído branco e criptografia. O DataSynthDecripter é uma ferramenta complementar que permite a recuperação dos dados originais a partir dos dados encriptados. Este documento detalha cada trecho do código, explicando suas funcionalidades e conceitos subjacentes, incluindo a transformação de Box-Muller para gerar ruído branco gaussiano.

2 Importações e Bibliotecas

O código utiliza várias bibliotecas para manipulação de dados, visualização e criptografia:

```
1 import pandas as pd
2 import numpy as np
3 import streamlit as st
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from PIL import Image
7 import base64
8 from io import BytesIO
9 from cryptography.fernet import Fernet
10 import hashlib
11 import random
12 import string
```

- pandas e numpy são usados para manipulação de dados.
- streamlit é utilizado para criar a interface web.
- matplotlib e seaborn são usados para visualização de dados.
- PIL é utilizado para manipulação de imagens.
- base64 e io são usados para codificação e manipulação de fluxos de bytes.
- cryptography.fernet é utilizado para criptografia simétrica.
- hashlib e string são usados para derivar chaves e gerar senhas.

3 Funções de Geração de Ruído Branco

3.1 Transformação de Box-Muller

A transformação de Box-Muller é uma técnica para gerar variáveis aleatórias com distribuição normal (gaussiana) a partir de variáveis aleatórias uniformemente distribuídas. A função gerar_ruido_branco implementa essa técnica:

```

1 def gerar_ruido_branco(tamanho):
2     U1 = np.random.uniform(0, 1, tamanho)
3     U2 = np.random.uniform(0, 1, tamanho)
4     Z0 = np.sqrt(-2 * np.log(U1)) * np.cos(2 * np.pi * U2)
5     return Z0

```

Explicação:

- U1 e U2 são variáveis aleatórias uniformemente distribuídas entre 0 e 1.
- Z0 é uma variável aleatória com distribuição normal gerada pela transformação de Box-Muller.
- A transformação utiliza funções logarítmica e trigonométricas para converter as variáveis uniformes em variáveis normais.

4 Funções de Criptografia

4.1 Derivar Chave a partir da Senha

A função `derivar_chave` gera uma chave de criptografia a partir de uma senha usando PBKDF2 com HMAC-SHA256:

```

1 def derivar_chave(senha):
2     kdf_salt = b'some_salt_' # Isso deve ser armazenado de
3     maneira segura
4     kdf_iterations = 100000
5     chave = base64.urlsafe_b64encode(hashlib.pbkdf2_hmac('sha256',
6     senha.encode(), kdf_salt, kdf_iterations, dklen=32))
7     return chave

```

Explicação:

- `senha` é a senha fornecida pelo usuário.
- `kdf_salt` é o salt usado na derivação da chave.
- `kdf_iterations` é o número de iterações para a função PBKDF2.
- `hashlib.pbkdf2_hmac` deriva uma chave segura de 256 bits.
- `base64.urlsafe_b64encode` codifica a chave derivada em um formato seguro para URLs.

4.2 Adicionar Ruído Branco e Encriptar Dados

A função `adicionar_ruido_branco_e_encriptar` adiciona ruído branco a uma coluna de dados e encripta os dados originais:

```

1 def adicionar_ruido_branco_e_encriptar(dados, coluna, chave,
2     ruido_nivel=1):
3     tamanho = len(dados)
4     Z0 = gerar_ruido_branco(tamanho)

```

```

4      dados[coluna] = dados[coluna].apply(lambda x: float(str(x).
      replace(',', ' ')))
5      ruido = Z0 * ruido_nivel * dados[coluna].std()
6      dados_sinteticos = dados.copy()
7      dados_sinteticos[coluna] += ruido
8
9      # Encriptar os dados originais
10     fernet = Fernet(chave)
11     dados['Encrypted_' + coluna] = dados[coluna].apply(lambda x:
      fernet.encrypt(str(x).encode()).decode())
12     return dados_sinteticos, dados

```

Explicação:

- dados é o DataFrame contendo os dados.
- coluna é a coluna onde o ruído será adicionado.
- chave é a chave de criptografia.
- ruido_nivel define a intensidade do ruído adicionado.
- dados_sinteticos contém os dados com ruído adicionado.
- Os dados originais são encriptados usando Fernet.

5 Funções de Anonimização

5.1 Anonimizar Dados Textuais com Ruído Branco

A função `anonimizar_texto_com_ruido` anonimiza dados textuais adicionando ruído:

```

1 def anonimizar_texto_com_ruido(dados, colunas_texto, chave):
2     dados_anonimizados = dados.copy()
3     fernet = Fernet(chave)
4     for coluna in colunas_texto:
5         dados_anonimizados['Encrypted_' + coluna] =
            dados_anonimizados[coluna].apply(lambda x: fernet.
            encrypt(str(x).encode()).decode())
6         dados_anonimizados[coluna] = gerar_imagem_ruido()
7     return dados_anonimizados

```

Explicação:

- dados é o DataFrame original.
- colunas_texto são as colunas textuais a serem anonimizadas.
- dados_anonimizados contém os dados anonimizados.
- Os valores textuais originais são encriptados e substituídos por imagens de ruído.

5.2 Gerar Imagem de Ruído Branco

A função `gerar_imagem_ruído` cria uma imagem de ruído branco:

```
1 def gerar_imagem_ruído(branco=True):
2     largura, altura = 50, 50
3     array_ruído = np.random.randint(0, 255, (altura, largura),
4         dtype=np.uint8)
5     img_ruído = Image.fromarray(array_ruído)
6     buffer = BytesIO()
7     img_ruído.save(buffer, format="PNG")
8     img_b64 = base64.b64encode(buffer.getvalue()).decode()
9     return f''
```

Explicação:

- `array_ruído` gera uma matriz de valores aleatórios.
- `Image.fromarray` cria uma imagem a partir da matriz de ruído.
- A imagem é codificada em base64 para ser exibida no HTML.

6 Funções de Interface do Streamlit

6.1 Ajustar Estilo CSS

A função `ajustar_estilo` aplica estilos CSS à interface:

```
1 def ajustar_estilo():
2     st.markdown(
3         """
4         <style>
5         .reportview-container .main .block-container {
6             max-width: 80%;
7             margin: auto;
8             padding: 2rem;
9         }
10        .dataframe {
11            width: 100%;
12            overflow-x: auto;
13            display: block;
14            white-space: nowrap;
15        }
16        </style>
17        """,
18        unsafe_allow_html=True
19    )
```

Explicação:

- Ajusta o layout da interface Streamlit para melhor visualização.

6.2 Formatar Colunas do DataFrame

A função `formatar_colunas` formata colunas específicas do DataFrame:

```
1 def formatar_colunas(dados):
2     dados['Codigo_Autorizacao'] = dados['Codigo_Autorizacao'].
      astype(int)
3     dados['Valor_Transacao'] = dados['Valor_Transacao'].apply(
      lambda x: float(str(x).replace(',', ' ')))
4     dados['Valor_Transacao'] = dados['Valor_Transacao'].apply(
      lambda x: f'{x:,.2f}'.replace(',', 'v').replace('.', ',').
      replace('v', '.'))
5     return dados
```

Explicação:

- Converte colunas para tipos específicos e formata valores para melhor exibição.

6.3 Gerar Senha Aleatória

A função `gerar_senha_aleatoria` gera uma senha aleatória:

```
1 def gerar_senha_aleatoria(tamanho=12):
2     caracteres = string.ascii_letters + string.digits
3     senha = ''.join(random.choice(caracteres) for i in range(
      tamanho))
4     return senha
```

Explicação:

- Gera uma senha aleatória de comprimento especificado, combinando letras e dígitos.

6.4 Decriptar Dados

A função `decriptar_dados` decripta os dados encriptados:

```
1 def decriptar_dados(dados, chave, colunas_texto):
2     fernet = Fernet(chave)
3     for coluna in dados.columns:
4         if coluna.startswith('Encrypted_'):
5             original_coluna = coluna.replace('Encrypted_', '')
6             if original_coluna in colunas_texto:
7                 try:
8                     dados[original_coluna] = dados[coluna].apply(
9                         lambda x: fernet.decrypt(x.encode()).decode
10                        ())
11                 except Exception as e:
12                     st.write(f"Erro ao decriptar a coluna {coluna
13                        }: {e}")
14             else:
15                 try:
16                     dados[original_coluna] = dados[coluna].apply(
17                         lambda x: float(fernet.decrypt(x.encode()).
18                        decode()))
```

```

14         except Exception as e:
15             st.write(f"Erro ao decriptar a coluna {coluna
16                     }: {e}")
17
18         # Remover colunas encriptadas ap s decripta o
19         colunas_encriptadas = [col for col in dados.columns if col.
20                               startswith('Encrypted_')]
21         dados.drop(columns=colunas_encriptadas, inplace=True)
22         return dados

```

Explicação:

- Decripta as colunas que foram encriptadas, convertendo-as de volta para seus valores originais.

7 Fluxo Principal do DataSynth

7.1 Interface do Streamlit

A interface principal do Streamlit é configurada para carregar um arquivo CSV, aplicar anonimização, adicionar ruído branco, encriptar os dados e permitir a decriptação dos dados encriptados:

```

1 ajustar_estilo()
2
3 st.title("DataSynth - Demonstra o de Anonimiza o de Dados
4     Financeiros")
5
6 st.write("Veja como os dados sens veis s o anonimizados para
7     proteger a privacidade dos clientes.")
8
9
10 # Inicializa a senha e chave usando o estado do Streamlit
11 if 'senha_aleatoria' not in st.session_state:
12     st.session_state.senha_aleatoria = None
13
14 if 'chave' not in st.session_state:
15     st.session_state.chave = None
16
17 if 'uploaded_file_buffer' not in st.session_state:
18     st.session_state.uploaded_file_buffer = None
19
20
21 # Bot o para gerar nova sess o
22 def gerar_nova_sessao():
23     st.session_state.senha_aleatoria = None
24     st.session_state.chave = None
25     st.session_state.uploaded_file_buffer = None
26     st.experimental_rerun()
27
28
29 # Carregar banco de dados de um arquivo CSV
30 uploaded_file = st.file_uploader("Escolha um arquivo CSV", type="
31     csv")
32
33 if uploaded_file is not None:
34     # Verifica se o arquivo carregado novo

```

```

27 uploaded_file_buffer = uploaded_file.getvalue()
28 if uploaded_file_buffer != st.session_state.
    uploaded_file_buffer:
29     # Limpa o estado atual
30     st.session_state.uploaded_file_buffer =
        uploaded_file_buffer
31
32     # Gera o da senha aleatoria ap s o upload do CSV
33     st.session_state.senha_aleatoria = gerar_senha_aleatoria()
34
35     # Gera o da chave secreta ap s o upload do CSV
36     st.session_state.chave = derivar_chave(st.session_state.
        senha_aleatoria)
37
38     # Limpa as tabelas e entradas de senha
39     st.session_state['banco_dados_transacoes'] = None
40     st.session_state['banco_dados_anonimizados'] = None
41     st.session_state['banco_dados_sinteticos'] = None
42     st.session_state['banco_dados_encryptados'] = None
43     st.session_state['senha_input'] = None
44
45 fernet = Fernet(st.session_state.chave)
46
47 st.write(f"Senha para decriptar: {st.session_state.
    senha_aleatoria}")
48
49 banco_dados_transacoes = pd.read_csv(BytesIO(st.session_state.
    uploaded_file_buffer))
50 banco_dados_transacoes['Valor_Transacao'] =
    banco_dados_transacoes['Valor_Transacao'].apply(lambda x:
        float(str(x).replace(',', ' ')))
51 banco_dados_transacoes = formatar_colunas(
    banco_dados_transacoes)
52
53 # Armazenar os dados no estado
54 st.session_state['banco_dados_transacoes'] =
    banco_dados_transacoes
55
56 # Adiciona o bot o "Gerar Nova Sess o" logo abaixo do upload
    do CSV
57 st.button("Gerar Nova Sess o", on_click=gerar_nova_sessao)
58
59 # Exibir dados originais
60 st.write("### Dados Originais")
61 st.write(st.session_state['banco_dados_transacoes'].head())
62
63 # Anonimizar dados textuais com ru do branco
64 colunas_texto = ['Nome_Cliente', 'Numero_Cartao', '
    Email_Cliente', 'Telefone_Cliente', 'Endereco_IP']
65 banco_dados_anonimizados = anonimizar_texto_com_ruido(st.
    session_state['banco_dados_transacoes'], colunas_texto, st.

```



```

        session_state.chave)
66 banco_dados_anonimizados['Valor_Transacao'] =
        banco_dados_anonimizados['Valor_Transacao'].apply(lambda x:
            float(str(x).replace(',','.')))
67 banco_dados_anonimizados = formatar_colunas(
        banco_dados_anonimizados)
68
69 # Armazenar os dados anonimizados no estado
70 st.session_state['banco_dados_anonimizados'] =
        banco_dados_anonimizados
71
72 # Aplicar ru do branco e encriptar os dados originais
73 ruido_nivel = st.slider("N vel de Ru do", min_value=0,
        max_value=1000, value=50, step=10)
74 banco_dados_sinteticos, banco_dados_encriptados =
        adicionar_ruido_branco_e_encriptar(st.session_state['
            banco_dados_anonimizados'], 'Valor_Transacao', st.
            session_state.chave, ruido_nivel)
75 banco_dados_sinteticos = formatar_colunas(
        banco_dados_sinteticos)
76
77 # Armazenar os dados sinteticos e encriptados no estado
78 st.session_state['banco_dados_sinteticos'] =
        banco_dados_sinteticos
79 st.session_state['banco_dados_encriptados'] =
        banco_dados_encriptados
80
81 # Remover colunas encriptadas dos dados sint ticos
82 colunas_encriptadas = [col for col in st.session_state['
            banco_dados_sinteticos'].columns if col.startswith('
            Encripted_')]
83 banco_dados_sinteticos_sem_encriptacao = st.session_state['
            banco_dados_sinteticos'].drop(columns=colunas_encriptadas)
84
85 # Exibir dados anonimizados com ru do branco aplicado com
        controle deslizante horizontal
86 st.write("### Dados Anonimizados com Ru do Branco Aplicado")
87 st.write('<div style="overflow-x: auto;">' +
        banco_dados_sinteticos_sem_encriptacao.head(5).to_html(
            escape=False, index=False) + '</div>', unsafe_allow_html=
            True)
88
89 # Bot o para exportar dados encriptados
90 csv_encriptado = banco_dados_encriptados.to_csv(index=False).
        encode('utf-8')
91 st.download_button(
92     label="Exportar Dados Encriptados",
93     data=csv_encriptado,
94     file_name='dados_encriptados.csv',
95     mime='text/csv'
96 )

```

```

97
98 # Verifica o de senha e decripta o de dados
99 senha_input = st.text_input("Insira a senha para decriptar os
    dados", type="password")
100 decriptar_btn = st.button("Decriptar Dados", key="
    decriptar_btn")
101 if decriptar_btn and senha_input == st.session_state.
    senha_aleatoria:
102     banco_dados_decriptados = decriptar_dados(st.session_state
        ['banco_dados_encryptados'], st.session_state.chave,
        colunas_texto)
103     st.write("### Dados Decriptados")
104     st.write('<div style="overflow-x: auto;">' +
        banco_dados_decriptados.head(5).to_html(escape=False,
            index=False) + '</div>', unsafe_allow_html=True)
105 elif decriptar_btn:
106     st.write("Senha incorreta. Por favor, tente novamente.")
107
108 # Visualizar os dados originais e sintéticos
109 fig, ax = plt.subplots(1, 2, figsize=(12, 6))
110 ax[0].hist(st.session_state['banco_dados_transacoes']['
    Valor_Transacao'], bins=30, alpha=0.6, color='blue', label=
    'Original')
111 ax[0].set_title('Dados Originais')
112 ax[0].set_xlabel('Valor Transação')
113 ax[0].set_ylabel('Frequência')
114
115 ax[1].hist(st.session_state['banco_dados_sinteticos']['
    Valor_Transacao'].apply(lambda x: float(str(x).replace('.',
        '' ).replace(',', '.'))), bins=30, alpha=0.6, color='green'
    , label='Sintético')
116 ax[1].set_title(f'Dados Sintéticos (Nível de Ruído = {
    ruido_nivel})')
117 ax[1].set_xlabel('Valor Transação')
118 ax[1].set_ylabel('Frequência')
119
120 st.pyplot(fig)
121
122 # Gerar e visualizar a PDF do ruído branco
123 ruido_branco = gerar_ruido_branco(1000)
124 fig, ax = plt.subplots()
125 sns.histplot(ruido_branco, bins=30, kde=True, ax=ax)
126 mu, std = np.mean(ruido_branco), np.std(ruido_branco)
127 xmin, xmax = ax.get_xlim()
128 x = np.linspace(xmin, xmax, 100)
129 p = (1/(np.sqrt(2*np.pi)*std)) * np.exp(-0.5*((x-mu)/std)**2)
130 ax.plot(x, p, 'k', linewidth=2)
131 ax.set_title("Histograma e Curva da Distribuição Normal")
132 ax.set_xlabel("Valor")
133 ax.set_ylabel("Densidade")
134 st.pyplot(fig)

```

```

135     st.write("Histograma e Curva da Distribuição Normal para o  

           Ru do Branco")
136 else:
137     st.write("Por favor, carregue um arquivo CSV para continuar.")

```

8 Fluxo Principal do DataSynthDecrypter

8.1 Interface do Streamlit

O DataSynthDecrypter é uma ferramenta complementar que permite a recuperação dos dados originais a partir dos dados encriptados. A interface principal do Streamlit é configurada para carregar um arquivo CSV encriptado, solicitar a senha de deciptação, e exibir os dados deciptados:

```

1 import pandas as pd
2 import streamlit as st
3 from cryptography.fernet import Fernet
4 import base64
5 import hashlib
6 from io import BytesIO
7
8 # Função para derivar chave a partir da senha
9 def derivar_chave(senha):
10     kdf_salt = b'some_salt_' # Isso deve ser armazenado de
11                               maneira segura
12     kdf_iterations = 100000
13     chave = base64.urlsafe_b64encode(hashlib.pbkdf2_hmac('sha256',
14                                                         senha.encode(), kdf_salt, kdf_iterations, dklen=32))
15     return chave
16
17 # Função para deciptar os dados
18 def deciptar_dados(dados, chave, colunas_texto):
19     fernet = Fernet(chave)
20     for coluna in dados.columns:
21         if coluna.startswith('Encrypted_'):
22             original_coluna = coluna.replace('Encrypted_', '')
23             if original_coluna in colunas_texto:
24                 try:
25                     dados[original_coluna] = dados[coluna].apply(
26                         lambda x: fernet.decrypt(x.encode()).decode()
27                     )
28                 except Exception as e:
29                     st.write(f"Erro ao deciptar a coluna {coluna}: {e}")
30             else:
31                 try:
32                     dados[original_coluna] = dados[coluna].apply(
33                         lambda x: float(fernet.decrypt(x.encode()).decode())
34                     )
35                 except Exception as e:

```

```

30         st.write(f"Erro ao decriptar a coluna {coluna
31                 }: {e}")
32     # Remover colunas encriptadas ap s decripta o
33     colunas_encriptadas = [col for col in dados.columns if col.
34                           startswith('Encrypted_')]
35     dados.drop(columns=colunas_encriptadas, inplace=True)
36     return dados
37
38 # Interface do Streamlit
39 st.title("DataSynthDecrypter - Desencripta o de Dados")
40
41 # Carregar banco de dados encriptado de um arquivo CSV
42 uploaded_file = st.file_uploader("Escolha um arquivo CSV
43 encriptado", type="csv")
44 if uploaded_file is not None:
45     banco_dados_encriptados = pd.read_csv(uploaded_file)
46
47     # Verifica o de senha e decripta o de dados
48     senha_input = st.text_input("Insira a senha para decriptar os
49 dados", type="password")
50     if st.button("Decryptar Dados"):
51         chave = derivar_chave(senha_input)
52         try:
53             banco_dados_decriptados = decriptar_dados(
54                 banco_dados_encriptados, chave, ['Nome_Cliente', '
55             Numero_Cartao', 'Email_Cliente', 'Telefone_Cliente'
56             , 'Endereco_IP', 'Valor_Transacao'])
57             st.write("### Dados Decriptados")
58             st.write(banco_dados_decriptados.head())
59         except Exception as e:
60             st.write(f"Erro na decripta o: {e}")

```

Explicação:

- **Carregar Dados Encriptados:** O arquivo CSV encriptado é carregado através do `st.file_uploader`.
- **Solicitar Senha:** A senha de decriptação é solicitada ao usuário através do `st.text_input`.
- **Derivar Chave:** A chave de decriptação é derivada da senha usando a função `derivar_chave`.
- **Decryptar Dados:** Os dados encriptados são decriptados usando a função `decriptar_dados`, e os dados decriptados são exibidos na interface.