

```

1 import torch
2 import os
3
4 # =====
5 # CONFIGURAÇÕES INICIAIS
6 # =====
7 script_dir = os.path.dirname(os.path.abspath(__file__))
8 FILE_PATH = os.path.join(script_dir, '..', 'machado.txt') # Sobe um nível e acha machado.txt
9
10 BLOCK_SIZE = 8 # Tamanho do contexto (janela de tempo)
11 BATCH_SIZE = 4 # Quantas sequências processar em paralelo
12 TRAIN_SPLIT = 0.8 # 80% para treino, 20% para validação
13
14 # =====
15 # 1. CARREGAMENTO DOS DADOS E VOCABULÁRIO
16 # =====
17 def load_data(file_path):
18     """
19         Lê o arquivo, cria o vocabulário e converte o
20         texto para inteiros.
21         Equivalente ao trecho em R que usa readr:::
22         read_file e stringr::str_split.
23     """
24     try:
25         with open(file_path, 'r', encoding='utf-8') as f:
26             text = f.read()
27             print(f"Arquivo '{file_path}' carregado com sucesso.")
28             # Equivalente ao stringr::str_sub(data, 100,
29             # 1000) apenas para visualização
30             print(f"Amostra do texto: {text[100:200]}...")

```

```

27 )
28     except FileNotFoundError:
29         print(f"Arquivo '{file_path}' não encontrado")
30         . Usando texto de exemplo.")
31         text = "Aqui está um texto de exemplo para"
32         simular o funcionamento do código. " * 500
33
34
35
36     chars = sorted(list(set(text)))
37     vocab_size = len(chars)
38     print(f"Tamanho do vocabulário: {vocab_size}")
39
40
41     # Mapeamentos (Encoder/Decoder)
42     # R: vocab <- structure(seq_along(vocab), names
43     # = vocab) (cria índices nomeados)
44     # Python: Criamos dicionários explícitos para
45     # conversão char <-> int
46
47
48     # Função para codificar: string -> lista de
49     # inteiros
50     encode = lambda s: [stoi[c] for c in s]
51     # Função para decodificar: lista de inteiros ->

```

```
51 string
52     decode = lambda l: ''.join([itos[i] for i in l])
53
54
# -----
#
55 # Criação do Tensor de Dados
56 # R: data <- stringr::str_split_1(data, "") ...
# depois acessa via índices
57 # Python: Convertemos o texto para um tensor
# LongTensor (inteiros)
58
# -----
#
59 data_tensor = torch.tensor(encode(text), dtype=
torch.long)
60 print(f"Total de tokens no dataset (caracteres): {len(data_tensor)}")
61 print(f"Primeiros 18 tokens (caracteres): '{decode(data_tensor[:18].tolist())}'")
62
63 return data_tensor, vocab_size, encode, decode
64
65
66 # =====
# =====
67 # 2. SEPARAÇÃO TREINO / VALIDAÇÃO
68 # =====
# =====
69 # Carrega os dados
70 data, vocab_size, encode, decode = load_data(
FILE_PATH)
71
72 # Define o ponto de corte (n no R refere-se ao length
# (data))
73 n = len(data)
74 n_train = int(n * TRAIN_SPLIT)
75
76 # Divide os dados em dois tensores distintos
77 train_data = data[:n_train]
78 val_data = data[n_train:]
```

```

79
80
81 # =====
82 # 3. FUNÇÃO GET_BATCH
83 # =====
84 def get_batch(split, block_size=BLOCK_SIZE,
85               batch_size=BATCH_SIZE):
86     """
87     Gera um lote (batch) de entradas (x) e alvos (y).
88     Args:
89         split (str): 'train' ou 'valid' (qualquer
90                     coisa diferente de train).
91         block_size (int): Tamanho do contexto.
92         batch_size (int): Número de sequências no
93                           batch.
94
95     Returns:
96         x (torch.Tensor): Tensor de entrada (
97             batch_size, block_size)
98         y (torch.Tensor): Tensor alvo (batch_size,
99                           block_size)
100
101    # Seleciona o dataset correto
102    # R: if (split == "train") { sample... } else {
103        sample... }
104    data_source = train_data if split == 'train'
105    else val_data
106
107    # Gera índices aleatórios para o início de cada
108    # bloco
109    # R: sample.int(0.8*(n-block_size), batch_size)
110    # Python: torch.randint gera inteiros aleatórios
111    # até high (len(data_source) - block_size)
112    ix = torch.randint(len(data_source) - block_size
113                      , (batch_size,))
114
115    # Constrói o tensor X empilhando os pedaços do
116    # texto

```

```

106      # R: x <- vocab[data[batch_ids]]; dim(x) <- dim(
107          batch_ids)
107      # Python: Para cada índice i em ix, pegamos o
108          slice data[i : i+block_size]
108      x = torch.stack([data_source[i: i + block_size]
109          for i in ix])
109
110      # Constrói o tensor Y (alvo) deslocado em 1
110          posição à direita
111      # R: y <- vocab[data[batch_ids + 1]]
112      y = torch.stack([data_source[i + 1: i +
112          block_size + 1] for i in ix])
113
114      # R: list(x = torch_tensor(x), y = torch_tensor(
114          y))
115      # Python: No PyTorch, x e y já são tensores aqui
115          se data_source for tensor.
116      # Mas garantimos que estejam no dispositivo
116          correto (CPU/GPU) se necessário no futuro.
117      return x, y
118
119
120 # =====
120 # =====
121 # TESTE DA IMPLEMENTAÇÃO (Main)
122 # =====
122 # =====
123 if __name__ == "__main__":
124     torch.manual_seed(1337) # Para
124         reproducibilidade, assim como set.seed no R
125
126     print("\n--- Testando get_batch('train') ---")
127     xb, yb = get_batch('train')
128
129     print("Shape de X:", xb.shape) # Deve ser [4, 8
129         ] (batch_size, block_size)
130     print("Shape de Y:", yb.shape)
131
132     print("\nConteúdo de X (Inputs):")
133     print(xb)
134

```

```
135     print("\nConteúdo de Y (Targets/Labels):")
136     print(yb)
137
138     print("\n--- Verificação Visual (Decodificação
139         ) ---")
140     # Vamos pegar a primeira linha do batch para
141     # conferir se Y é o deslocamento de X
142     primeira_linha_x = xb[0].tolist()
143     primeira_linha_y = yb[0].tolist()
144
145     print(f"X[0] decodificado: '{decode(
146         primeira_linha_x)}'")
147     print(f"Y[0] decodificado: '{decode(
148         primeira_linha_y)}'")
149
150     print("Nota: Y deve ser X deslocado uma letra
151         para frente (next token prediction).")
```