



ISEL
INSTITUTO SUPERIOR DE
ENGENHARIA DE LISBOA

Projeto 1

Ambiente Virtuais Interativos e Inteligentes

Trabalho realizador por:

Duarte Valente | A47657

João Valido | A51090

Docente:

Eng. Arnaldo Abrantes

Curso: MEIM

2023

Índice

1. Introdução	3
a) Descrição	3
b) Âmbito	3
c) Documentação	3
A. Mãos na massa	4
1. Modelos baseados em bigramas	4
a) Contagens no conjunto de treino	4
b) Otimização do critério de ML	6
2. Modelo MLP	9
B. Conclusão	14

Índice de ilustrações

Figura 1 - Criação de bigramas	4
Figura 2 - Criação de variáveis	4
Figura 3 - Contagem de bigramas	4
Figura 4 - Código para gerar palavras	5
Figura 5 - Palavras obtidas	5
Figura 6 - Cálculo do loss	5
Figura 7 - Demonstração gráfica da aprendizagem do modelo	6
Figura 8 - Criação de bigramas	6
Figura 9 - Criação de variáveis	7
Figura 10 - Criar modelo de treino	7
Figura 11 - Treino do modelo	8
Figura 12 - Demonstração gráfica da aprendizagem do modelo	8
Figura 13 - criação de variáveis	9
Figura 14 - Definir datasets para treino, validação e teste	9
Figura 15 - Definição de parâmetros e weights e biases	10
Figura 16 - Treino do modelo	10
Figura 17 - Demonstração gráfica da aprendizagem do modelo	11
Figura 18 - Avaliação dos datasets	12
Figura 19 - Palavras geradas	13

1.Introdução

a) Descrição

Este trabalho tinha como objetivo introduzir o tema da IA Generativa, isto é, ferramentas com capacidade para produzir conteúdos novos, como textos, imagens, vídeos, etc., a partir de modelos de ML, treinados sobre um grande conjunto de dados. Tem ainda como objetivo refletir acerca do impacto destas ferramentas de IA na sociedade, nomeadamente nas indústrias de entretenimento e jogos.

b) Âmbito

Este relatório insere-se no âmbito da realização da ficha laboratorial unidade curricular de Ambientes Virtuais Interativos e Inteligentes, do Mestrado em Engenharia Informática e Multimédia do DEETC do ISEL.

c) Documentação

Documentação de apoio à unidade curricular de Ambientes Virtuais Interativos e Inteligentes.

A. Mãos na massa

1. Modelos baseados em bigramas

a) Contagens no conjunto de treino

- criar lista de todos os bigramas encontrados no conjunto de nomes e as respectivas vezes que ocorrem

```
1 # criar lista de todos os bigramas encontrados no conjunto de nomes e as respectivas vezes que ocorrem
2 b = {}
3 for w in words:
4     w = ['<s>'] + list(w) + ['<e>']
5     for ch1, ch2 in zip(w, w[1:]):
6         bigram = (ch1, ch2)
7         b[bigram] = b.get(bigram, 0) + 1
```

Figura 1 - Criação de bigramas

- criar a lista de todos os caracteres "chars", fazer as correspondências para inteiros "stoi", fazer a correspondência para string "itos"

```
1 # criar a lista de todos os caracteres "chars",
2 # fazer as correspondência para inteiros "stoi", fazer a correspondência para string "itos"
3 chars = sorted(list(set(''.join(words))))
4 chars = ['.' + chars
5 stoi = {s:i for i,s in enumerate(chars)}
6 itos = {i:s for i,s in enumerate(chars)}
```

Figura 2 - Criação de variáveis

- para cada palavra da lista de nomes, preencher a matriz N com o número de vezes que cada bigrama aparece

```
1 # para cada palavra da lista de nomes, preencher a matriz N com o numero de vezes que cada bigrama aparece
2 for w in words:
3     w = ['.' + list(w) + ['.']
4     for ch1, ch2 in zip(w, w[1:]):
5         ix1 = stoi[ch1]
6         ix2 = stoi[ch2]
7         N[ix1, ix2] += 1
```

Figura 3 - Contagem de bigramas

- **gerar 5 palavras com base na probabilidade da próxima letra a ser gerada**

```
1 # gerar 5 palavras com base na probabilidade da proxima letra a ser gerada
2 g = torch.Generator().manual_seed(123789)
3 for i in range(5):
4     out = []
5     ix = stoi['.']
6     while True:
7         pp = p[ix]
8         ix = torch.multinomial(pp, num_samples=1, replacement=True, generator=g).item()
9         out.append(itos[ix])
10    if ix == 0:
11        break
12    print(''.join(out))
```

Figura 4 - Código para gerar palavras

```
nariel.
sa.
sy.
a.
welesty.
```

Figura 5 - Palavras obtidas

```
1 lossi = []
2 nll = 0
3 n = 0
4 for w in words:
5     w = ['.'] + list(w) + ['.']
6     for ch1, ch2 in zip(w, w[1:]):
7         ix1 = stoi[ch1]
8         ix2 = stoi[ch2]
9         logprob = torch.log(p[ix1, ix2])
10        nll -= logprob
11        n += 1
12    lossi.append(nll/n)
13 print(nll/n)
```

Figura 6 - Cálculo do loss

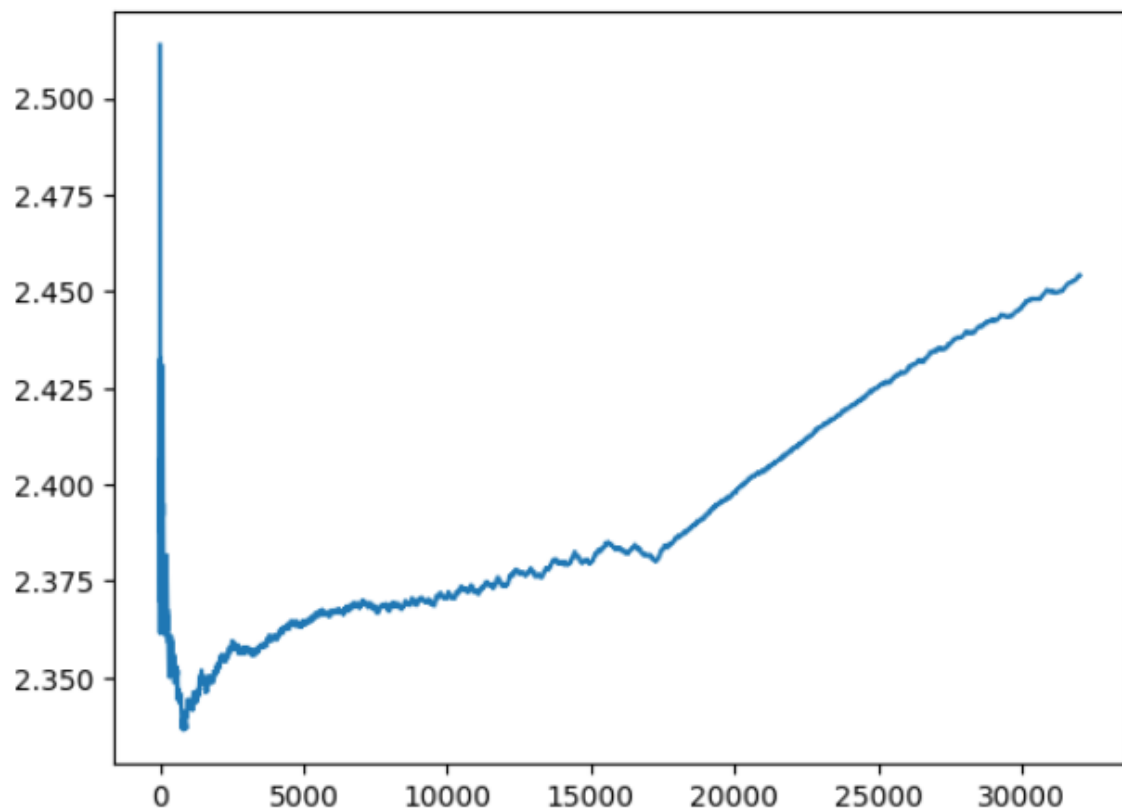


Figura 7 - Demonstração gráfica da aprendizagem do modelo

b) Otimização do critério de ML

- **criar lista de todos os bigramas encontrados no conjunto de nomes e as respetivas vezes que ocorrem**

```
1 # criar lista de todos os bigramas encontrados no conjunto de nomes e as respetivas vezes que ocorrem
2 b = {}
3 for w in words:
4     w = ['<s>'] + list(w) + ['<e>']
5     for ch1, ch2 in zip(w, w[1:]):
6         bigram = (ch1, ch2)
7         b[bigram] = b.get(bigram, 0) + 1
```

Figura 8 - Criação de bigramas

- **criar a lista de todos os caracteres "chars", fazer as correspondências para inteiros "stoi", fazer a correspondência para string "itos"**

```
1 # criar a lista de todos os caracteres "chars",
2 # fazer as correspondência para inteiros "stoi", fazer a correspondência para string "itos"
3 chars = sorted(list(set(''.join(words))))
4 chars = ['.'] + chars
5 stoi = {s:i for i,s in enumerate(chars)}
6 itos = {i:s for i,s in enumerate(chars)}
```

Figura 9 - Criação de variáveis

- **criar o modelo de treino**

```
1 # criar o modelo de treino
2 xs, ys =[], []
3 for w in words:
4     chs = ['.' ] + list(w) + ['.']
5     for ch1, ch2 in zip(chs, chs[1:]):
6         ix1 = stoi[ch1]
7         ix2 = stoi[ch2]
8         xs.append(ix1)
9         ys.append(ix2)
10 xs = torch.tensor(xs)
11 ys = torch.tensor(ys)
12 num = xs.shape[0]
```

Figura 10 - Criar modelo de treino

- **Treino do modelo**

```
1 # treino do modelo
2
3 for k in range(1000):
4     xenc = F.one_hot(xs, num_classes=27).float()
5     logits = xenc @ W
6     counts = logits.exp()
7     probs = counts / counts.sum(1, keepdim = True)
8
9     loss = -probs[torch.arange(xs.shape[0]), ys].log().mean()
10    lossi.append(loss.item())
11    W.grad = None
12    loss.backward()
13
14    W.data += -10*W.grad
15
16 print (loss)
```

Figura 11 - Treino do modelo

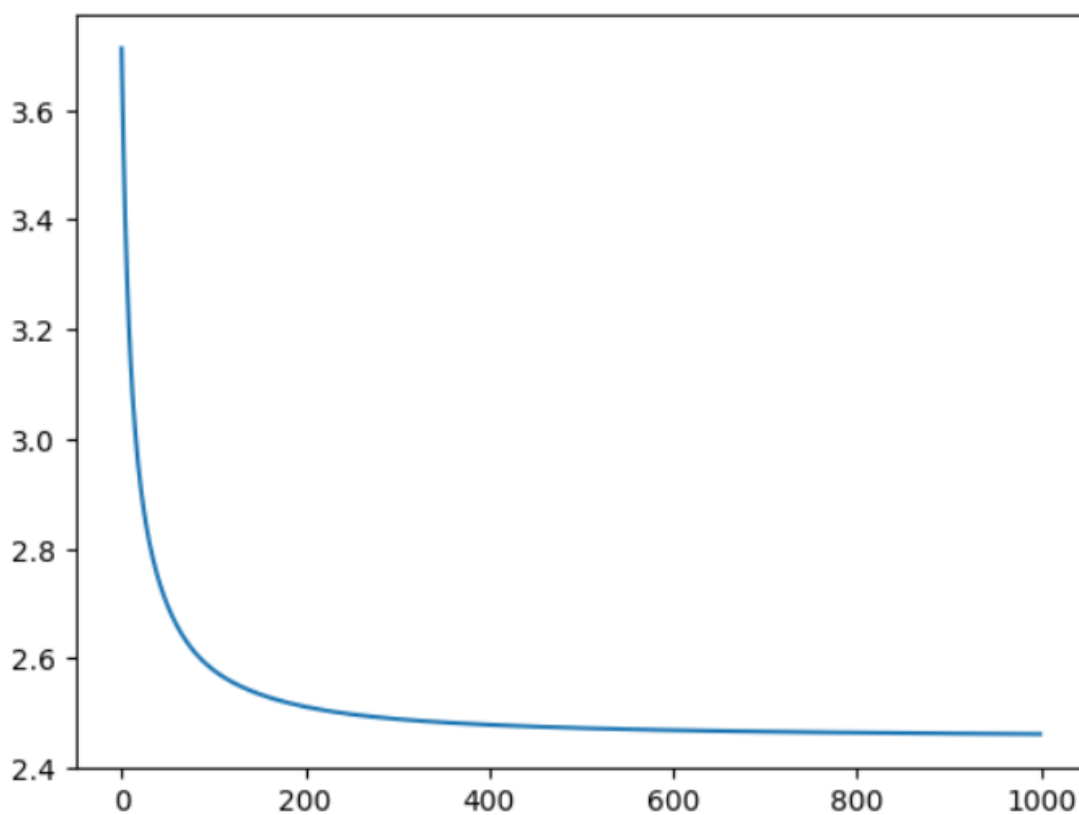


Figura 12 - Demonstração gráfica da aprendizagem do modelo

2. Modelo MLP

- criar a lista de todos os caracteres "chars", fazer as correspondências para inteiros "stoi", fazer a correspondência para string "itos"

```
1 # criar a lista de todos os caracteres "chars",
2 # fazer as correspondência para inteiros "stoi", fazer a correspondência para string "itos"
3 chars = sorted(list(set(''.join(words))))
4 chars = ['.' ] + chars
5 stoi = {s:i for i,s in enumerate(chars)}
6 itos = {i:s for s,i in stoi.items()}
```

Figura 13 - criação de variáveis

- Definir dataset para treino, validação e teste, 80%, 10%, 10% respetivamente

```
5 block_size = 3 # context length: how many characters do we take to predict the next one?
6
7 def build_dataset(words):
8     X, Y = [], []
9     for w in words:
10
11         #print(w)
12         context = [0] * block_size
13         for ch in w + '.':
14             ix = stoi[ch]
15             X.append(context)
16             Y.append(ix)
17             #print(''.join(itos[i] for i in context), '--->', itos[ix])
18             context = context[1:] + [ix] # crop and append
19
20 X = torch.tensor(X)
21 Y = torch.tensor(Y)
22 print(X.shape, Y.shape)
23 return X, Y
24
25 random.seed(42)
26 random.shuffle(words)
27 n1 = int(0.8*len(words))
28 n2 = int(0.9*len(words))
29
30 Xtr, Ytr = build_dataset(words[:n1])
31 Xdev, Ydev = build_dataset(words[n1:n2])
32 Xte, Yte = build_dataset(words[n2:])
```

Figura 14 - Definir datasets para treino, validação e teste

```
1 g = torch.Generator().manual_seed(2461359) # for reproducibility
2 C = torch.randn((27, 3), generator=g) # embedding space
3 W1 = torch.randn((9, 1000), generator=g)
4 b1 = torch.randn(1000, generator=g)
5 W2 = torch.randn((1000, 27), generator=g)
6 b2 = torch.randn(27, generator=g)
7 parameters = [C, W1, b1, W2, b2]
8 for p in parameters:
9     p.requires_grad = True
```

Figura 15 - Definição de parâmetros e weights e biases

```
1 for i in range(10000):
2     # minibatches
3     ix = torch.randint(Xtr.shape[0], (32,))
4
5     # forward pass
6     emb = C[Xtr[ix]]
7     h = torch.tanh(emb.view(-1, 9) @ W1 + b1)
8     logits = h @ W2 + b2
9     loss = F.cross_entropy(logits, Ytr[ix])
10    #print(i, loss.item())
11    lossi.append(loss.item())
12
13    # backward pass
14    for p in parameters:
15        p.grad = None
16    loss.backward()
17
18    # update
19    for p in parameters:
20        p.data += -0.1 * p.grad
21
22    print(loss.item())
```

1.9709446430206299

Figura 16 - Treino do modelo

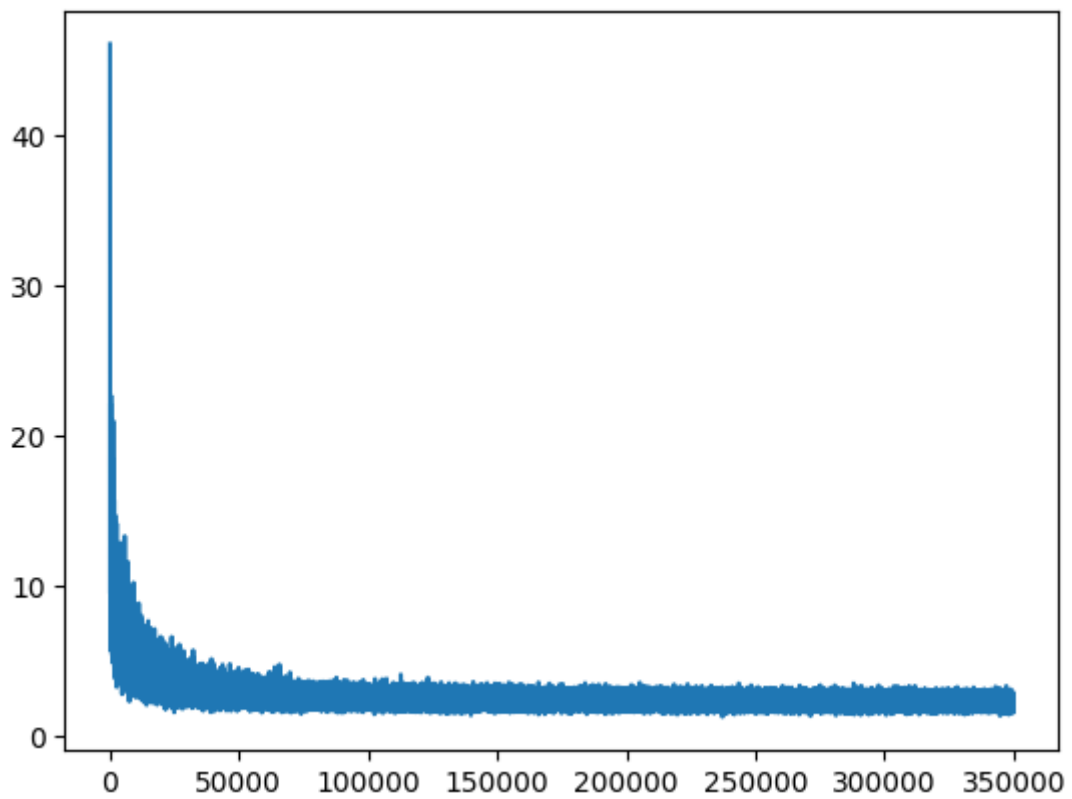


Figura 17 - Demonstração gráfica da aprendizagem do modelo

- **Através das avaliações dos datasets podemos verificar que o resultado da validação é um pouco maior que o resultado do treino, logo o modelo não se encontra em overfit**

```
1 # Evaluate in the training set
2 emb = C[Xtr]
3 h = torch.tanh(emb.view(-1, 9) @ W1 + b1)
4 logits = h @ W2 + b2
5 loss = F.cross_entropy(logits, Ytr)
6 loss

tensor(2.2858, grad_fn=<NllLossBackward0>)

1 # Evaluate in the validation set (is the model overfitting?)
2 emb = C[Xdev]
3 h = torch.tanh(emb.view(-1, 9) @ W1 + b1)
4 logits = h @ W2 + b2
5 loss = F.cross_entropy(logits, Ydev)
6 loss

tensor(2.3150, grad_fn=<NllLossBackward0>)

1 # Evaluate in the test set
2 emb = C[Xte]
3 h = torch.tanh(emb.view(-1, 9) @ W1 + b1)
4 logits = h @ W2 + b2
5 loss = F.cross_entropy(logits, Yte)
6 loss

tensor(2.3329, grad_fn=<NllLossBackward0>)
```

Figura 18 - Avaliação dos datasets

```
1 # Sampling from the model
2 g = torch.Generator().manual_seed(2461359)
3
4 for _ in range(20):
5
6     out = []
7     context = [0] * block_size # initialize with all ...
8     while True:
9         emb = C[torch.tensor([context])] # (1,block_size,d)
10        h = torch.tanh(emb.view(1, -1) @ W1 + b1)
11        logits = h @ W2 + b2
12        probs = F.softmax(logits, dim=1)
13        ix = torch.multinomial(probs, num_samples=1, generator=g).item()
14        context = context[1:] + [ix]
15        out.append(ix)
16        if ix == 0:
17            break
18
19    print(''.join(itos[i] for i in out))
```

ale.
maraus.
annase.
even.
elyn.
eve.
dite.
dephillanassimaysany.
ryxa.
ia.
lynn.
jayron.
meilyn.
aulaj.
san.
y.
sana.
eveleigany.
zaryas.
saur.

Figura 19 - Palavras geradas

B. Conclusão

Em conclusão, a comparação de modelos baseados em bigramas através da realização de contagens no conjunto de treino e da otimização do critério de ML usando o algoritmo de descida de gradiente e o modelo MLP com camada de "Embedding" foi um estudo valioso e esclarecedor. Os resultados mostraram que o modelo MLP com camada de "Embedding" apresentou melhor desempenho na tarefa de predição de palavras em um corpus de texto, indicando a importância de incluir uma camada de representação de palavras em modelos de processamento de linguagem natural. No entanto, é importante ressaltar que a escolha do modelo ideal dependerá do contexto e dos objetivos específicos da análise. Portanto, estudos adicionais são necessários para aprimorar ainda mais os modelos e suas aplicações em diferentes cenários, uma vez que também verificámos que um dataset mais volumoso não indica diretamente uma melhor aprendizagem por parte do modelo.