

University of Central Missouri
Department of Computer Science & Cybersecurity

CS5760 Natural Language Processing

Fall 2025

Homework 4.

Student name: Ajith.Bonda

Submission Requirements:

- Once finished your assignment push your source code to your repo (GitHub) and explain the work through the ReadMe file properly. Make sure you add your student info in the ReadMe file.
- Submit your GitHub link on the Bright Space.
- Comment your code appropriately **IMPORTANT**.
- Any submission after provided deadline is considered as a late submission.

GIT Link: [https://github.com/Ajith0028/NLP Assignment 4](https://github.com/Ajith0028/NLP_Assignment_4)

Part I. Short Answer

1) RNN Families & Use-Cases (Many-to-X)

a) Map each task to the most suitable RNN I/O pattern and explain why (1–2 lines each):

- next-word prediction

- sentiment of a sentence

- NER

- machine translation

(Choose from: one-to-many, many-to-one, many-to-many aligned, many-to-many unaligned.)

Ans:

Next-word prediction: many-to-many (aligned) - each input token produces a next-token prediction aligned in time.

Sentiment of a sentence: many-to-one - full sequence processed to output a single sentiment label.

NER: many-to-many (aligned) - every input word requires a corresponding output tag.

Machine translation: many-to-many (unaligned) - input and output sequences differ in length, requiring an encoder–decoder setup.

b) In one sentence, explain how “unrolling” over time enables BPTT and weight sharing.

Ans: Unrolling converts the recurrent loop into a time-step chain, enabling BPTT and allowing the same weights to be reused at every step.

c) Give one advantage and one limitation of weight sharing across time in RNNs.

Ans: **Advantage:** Fewer parameters and better generalization by learning time-invariant patterns.

Limitation: It cannot model time-specific behavior since all steps use the same shared weights

2) Vanishing Gradients & Remedies

a) Describe the vanishing gradient problem in RNNs and how it affects long-range dependencies.

Ans: Vanishing gradients occur when gradients shrink exponentially during backpropagation through many time steps, causing the RNN to fail to learn long-range dependencies because early time-step weights receive almost no gradient signal.

b) List two architectural solutions and briefly how each helps gradient flow.

Ans: **LSTM:** Uses gates and a cell state that provides a constant error path, preventing gradients from shrinking over long sequences.

GRU: Uses update/reset gates that regulate information flow, helping preserve gradients across time.

c) Give one training technique (not an architecture change) that can mitigate the issue and why.

Ans: **Gradient clipping:** Prevents exploding/unstable updates, stabilizing backpropagation and helping gradients stay in a usable range over long sequences.

3) LSTM Gates & Cell State

a) Explain the roles of the forget, input, and output gates. For each, name its activation and purpose.

Ans: **Forget gate** — σ (sigmoid): decides what portion of the previous cell state to erase.

Input gate — σ (sigmoid) + \tanh : decides which new information to add to the cell state.

Output gate — σ (sigmoid) + \tanh : controls how much of the cell state is revealed as the hidden state/output.

b) Why is the LSTM cell state often described as providing a “linear path” for gradients?

Ans: The LSTM cell state updates with element-wise additions and scaling instead of repeated multiplications, allowing gradients to pass through many time steps without vanishing.

c) In one or two sentences, contrast “what to remember” vs. “what to expose” in LSTMs.

Ans: “What to remember”- is controlled by the forget and input gates, which decide how the cell state is updated.

“What to expose” - is controlled by the output gate, which determines how much of the internal memory becomes visible as the hidden state.

4) Self-Attention

a) Define Query (Q), Key (K), and Value (V) in the context of self-attention.

Ans: **Query (Q):** a vector representing *what a token is looking for* in other tokens.

Key (K): a vector representing *what each token offers* or how it should be matched against queries.

Value (V): a vector containing the *actual information/content* to be aggregated after attention scores are computed.

b) Write the formula for dot-product attention.

Ans: $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k})V$

c) Why do we divide by $\sqrt{d_k}$?

Ans: To prevent the dot products from growing too large as dimensionality increases, which keeps the softmax stable and avoids extremely small gradients.

5) Multi-Head Attention & Residual Connections

a) Why do Transformers use multi-head attention instead of single-head attention?

Ans: Multi-head attention lets the model attend to different types of relationships (positions, syntax, meaning) in parallel, giving richer representations than a single-head can capture.

b) What is the purpose of Add & Norm (Residual + LayerNorm)? Explain two benefits.

Ans: **Residual** connections help preserve gradients and prevent degradation when stacking many layers.

LayerNorm stabilizes activations and speeds up training by normalizing each layer's outputs.

c) Describe one example of linguistic relation that different heads might capture (e.g., coreference, syntax).

Ans: One head might track coreference (e.g., “John... he”), while another tracks syntactic dependencies like subject–verb or adjective–noun relationships.

6) Encoder–Decoder with Masked Attention

a) Why does the decoder use masked self-attention? What problem does it prevent?

Ans: To prevent the decoder from looking at future tokens, ensuring it only predicts based on past outputs and avoiding information leakage.

b) What is the difference between encoder self-attention and encoder–decoder cross-attention?

Ans: **Encoder self-attention:** tokens attend only to other input tokens.

Encoder–decoder cross-attention: decoder queries attend to encoder outputs to use source-sentence information.

c) During inference (no teacher forcing), how does the model generate tokens step by step?

Ans: The model generates one token at a time, feeds it back into the decoder, and repeats the process until an end-of-sentence token is produced.

Part II: Programming

Q1. Character-Level RNN Language Model (“hello” toy & beyond)

Goal: Train a tiny character-level RNN to predict the next character given previous characters.

Data (toy to start):

- Start with a small toy corpus you create (e.g., several “hello...”, “help...”, short words/sentences).
- Then expand to a short plain-text file of ~50–200 KB (any public-domain text of your choice).

Model:

- Embedding → RNN (Vanilla RNN or GRU or LSTM) → Linear → Softmax over characters.

- Hidden size 64–256; sequence length 50–100; batch size 64; train 5–20 epochs.

Train:

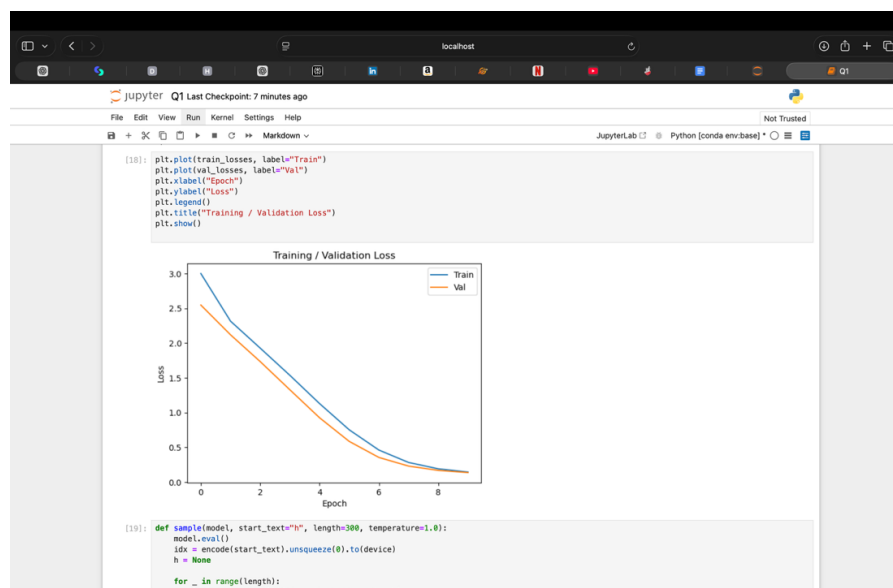
- Teacher forcing (use the true previous char as input during training).
- Cross-entropy loss; Adam optimizer.

Report:

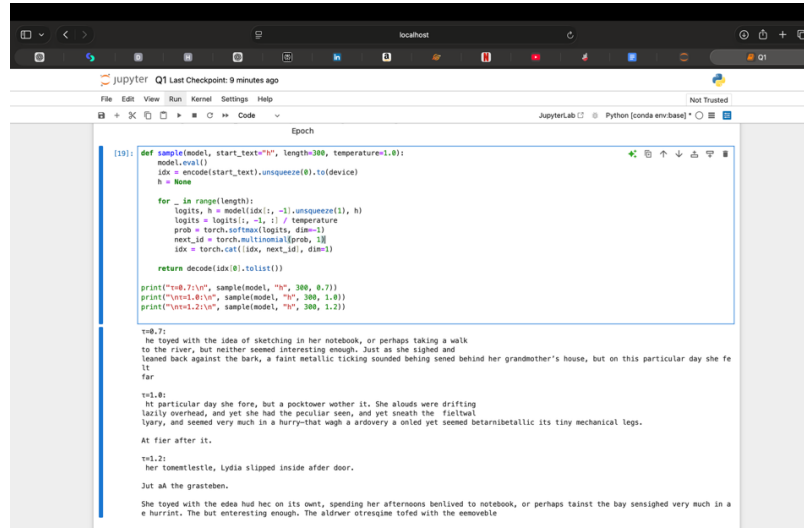
1. Training/validation loss curves.
2. Sample 3 temperature-controlled generations (e.g., $\tau = 0.7, 1.0, 1.2$) for 200–400 chars each.
3. A 3–5 sentence reflection: what changes when you vary sequence length, hidden size, and temperature?
(Connect to slides: embedding, sampling loop, teacher forcing, tradeoffs)

Ans:

1.



2.



The screenshot shows a JupyterLab window with a code editor and an output area. The code defines a function `sample(model, start_text="w", length=300, temperature=1.0)` that generates text using an RNN model. The output shows the generated text for three different temperatures: 0.7, 1.0, and 1.2.

```
[In]: def sample(model, start_text="w", length=300, temperature=1.0):
      model.eval()
      idx = encode(start_text).unsqueeze(0).to(device)
      h = None

      for _ in range(length):
          logits, h = model(idx.unsqueeze(0), h)
          logits = logits[:, -1, :] / temperature
          prob = torch.softmax(logits, dim=-1)
          next_id = torch.multinomial(prob, 1)
          idx = torch.cat([idx, next_id], dim=-1)

      return decode(idx[0].tolist())

      print("\n0.7:\n", sample(model, "w", 300, 0.7))
      print("\n1.0:\n", sample(model, "w", 300, 1.0))
      print("\n1.2:\n", sample(model, "w", 300, 1.2))

      "\n0.7:
      He toyed with the idea of sketching in her notebook, or perhaps taking a walk
      to the river, but neither seemed interesting enough. Just as she sighed and
      leaned back against the bark, a faint metallic ticking sounded behind her grandmother's house, but on this particular day she
      it
      far

      "\n1.0:
      He particular day she fore, but a pocketwatch mother it. She alouds were drifting
      lazily overhead, and yet she had the peculiar seen, and yet smelt the fieldwal
      lypay, and seemed very much in a hurry-that wagh a ardoverly a oned yet seemed betarnibetallic its tiny mechanical legs.

      At tier after it.

      "\n1.2:
      her tomentlestle, Lydia slipped inside afder door.

      Jut aa the grasteben.

      She toyed with the idea had hac on its omt, spending her afternoons belnived to notebook, or perhaps tainst the bay sentsighed very much in a
      e hurrint. The but interesting enough. The aldrwer atresqine tofed with the emmoveble
```

3. Increasing the sequence length allows the model to capture longer-range dependencies, improving coherence but also increasing training time and difficulty. A larger hidden size usually improves the model's capacity to learn richer patterns, but it may overfit and requires more memory. During sampling, temperature strongly affects creativity: low τ (≈ 0.7) produces more repetitive and conservative text, $\tau = 1.0$ gives balanced diversity, and high τ ($1.2+$) produces creative but often incoherent sequences. Teacher forcing helps the model converge faster by providing the correct previous character, but it can make sampling more brittle when the model must rely on its own predictions. Overall, the RNN learns meaningful character-level structure, but its generation quality is sensitive to these hyperparameters.

Q2. Mini Transformer Encoder for Sentences

Task: Build a mini Transformer Encoder (NOT full decoder) to process a batch of sentences.

Steps:

1. Use a small dataset (e.g., 10 short sentences of your choice).
2. Tokenize and embed the text.
3. Add sinusoidal positional encoding.
4. Implement:
 - o Self-attention layer
 - o Multi-head attention (2 or 4 heads)
 - o Feed-forward layer

- Add & Norm
5. Show:
- Input tokens
 - Final contextual embeddings
 - Attention heatmap between words (visual or printed)

Ans:

1.

```
[5]: # simple whitespace tokenizer
tokens = s.split() for s in sentences

# build vocabulary
vocab = sorted(w for sent in tokens for w in sent)
stoi = {w:i for i,w in enumerate(vocab)} # start indexing from 1
stoi["<pad>"] = 0
itos = {i:w for w,i in stoi.items()}

# encode
encoded = [[stoi[w] for w in sent] for sent in tokens]
max_len = max(len(x) for x in encoded)

# pad
padded = [seq + [0]*(max_len - len(seq)) for seq in encoded]
padded = torch.tensor(padded)
print("Input token indices:\n", padded)

Input token indices:
tensor([[30, 5, 27, 22, 30, 19],
        [30, 8, 6, 30, 5, 0],
        [1, 18, 15, 33, 13, 0],
        [28, 16, 31, 24, 18, 0],
        [30, 29, 15, 3, 0, 0],
        [2, 9, 14, 30, 29, 0],
        [12, 20, 1, 4, 0, 0],
        [30, 7, 15, 25, 0, 0],
        [30, 17, 23, 13, 0, 0],
        [30, 34, 15, 21, 32, 0]])

[7]: d_model = 32
vocab_size = len(stoi)

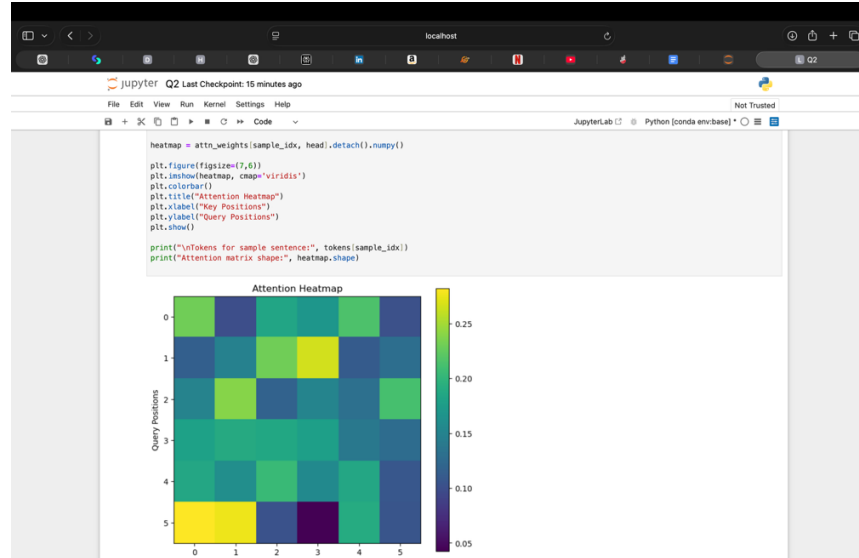
embed = nn.Embedding(vocab_size, d_model)

def sinusoidal_pos_encoding(seq_len, d_model):
    pos = torch.arange(seq_len).unsqueeze(1)
    i = torch.arange(d_model).unsqueeze(0)
    angle_rates = 1 / (10000 ** (2 * (i // 2) / d_model))
```

2.

```
Final contextual embeddings:
tensor([[[ 1.4665, 1.5208, 1.8348, ..., 1.4069, -0.3818, 2.1276],
         [ 0.7891, 0.2080, 0.5484, ..., 0.0787, -0.0485, 0.9725],
         [-0.6432, -0.0647, 0.4856, ..., -0.4359, -0.9229, -0.8751],
         [ 0.4779, -2.3648, 0.4086, ..., 1.3602, -0.4269, 0.2822],
         [ 1.0582, 0.1632, 2.0237, ..., 1.6286, -0.3893, 2.4377],
         [-1.3682, -1.3744, -0.5848, ..., -0.3538, -0.6469, 0.9178]],
        [[ 1.6228, 1.4684, 0.8857, ..., 1.3611, -0.1385, 2.2473],
         [-2.7717, -0.8558, 0.8778, ..., -0.2348, -0.3436, 0.8813],
         [-1.8228, -1.8464, 1.3176, ..., -0.4827, -0.5454, -1.1545],
         [ 1.6494, -0.2878, 1.9226, ..., 1.5898, -0.3899, 2.4131],
         [-0.8479, -1.1818, 0.9625, ..., 0.2818, -0.1366, 1.2863],
         [-2.4063, -1.5921, 0.4194, ..., -0.0734, 1.1288, 0.6123]],
        [[ 0.8412, 0.2351, 0.1679, ..., 2.2695, -0.7823, -0.1336],
         [ 0.1419, -0.1842, -1.5794, ..., 0.7886, -0.0408, 1.5188],
         [-0.2861, 0.1184, -1.2426, ..., 1.7175, -0.1879, -0.7961],
         [-0.2213, -1.6594, -0.9982, ..., 1.6427, -2.5805, -0.2371],
         [ 0.1469, -0.9189, -0.1968, ..., 0.9644, -0.7839, 0.5873],
         [-2.4688, -1.5825, 0.2224, ..., 0.8077, 1.2681, 0.1788]],
        ...,
        [[ 1.7289, 1.3645, 0.7412, ..., 1.5791, -0.3231, 2.2888],
         [ 1.5825, -1.2122, -1.1329, ..., 1.6638, -1.8488, -0.8391],
         [-0.2587, -0.3155, -1.2477, ..., 1.6785, -0.3858, -0.5816],
         [-0.9856, -1.3344, 0.3832, ..., -0.2518, -2.8863, -0.6998],
         [-2.8728, -2.3844, 0.6736, ..., 0.8443, 0.3943, -0.4719],
         [-2.3166, -1.7555, 0.3144, ..., 0.8118, 1.0724, 0.5221]],
        [[-1.7835, 0.8448, -2.4856, ..., -0.2138, 0.8655, 0.1285],
         [ 0.9135, -0.9055, -1.1823, ..., 1.8409, -3.0884, -0.6361],
         [-1.8767, -1.2194, 1.2423, ..., 0.8686, -2.9888, 0.3498],
         [-0.9853, -2.7969, -0.8119, ..., 0.8536, -1.9172, -0.6188],
         [-1.9944, -2.4444, 0.6492, ..., 0.8578, 0.9494, 0.2427],
         [-2.2393, -1.7989, 0.2741, ..., 0.8325, 1.8143, 0.2855]],
        [[ 1.6888, 1.4333, 0.7658, ..., 1.6165, -0.4261, 2.8868],
         [ 0.2837, -0.9768, 0.7242, ..., 0.8674, -1.1422, 0.3928],
         [-0.3886, -0.1146, -1.2571, ..., 1.6699, -0.3727, -0.5743],
         [-0.8878, -2.4163, -0.2195, ..., 1.1582, -0.8211, -0.1317],
         [-1.7431, -1.3937, -0.1387, ..., 0.8136, -0.6533, 1.5648],
         [-2.3765, -1.6221, 0.2575, ..., 0.8333, 1.8383, 0.4777]]])
```

3.



Q3. Implement Scaled Dot-Product Attention

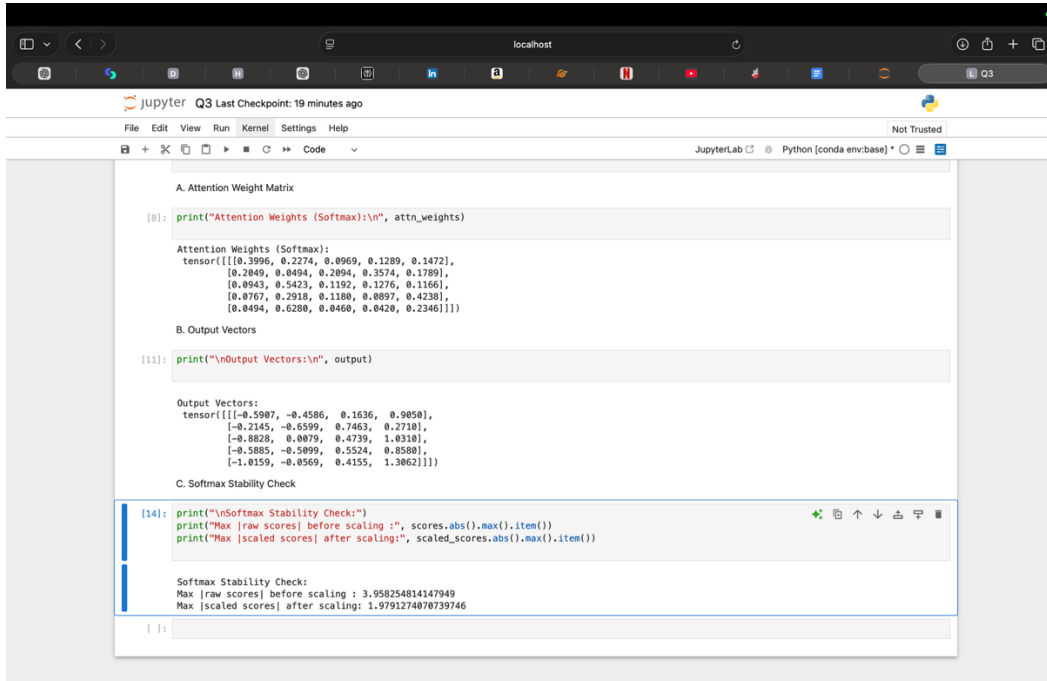
Goal: Implement the attention function from your slides:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Requirements:

- Write a function in PyTorch or TensorFlow to compute attention.
- Test it using random Q, K, V inputs.
- Print:
 - Attention weight matrix
 - Output vectors
 - Softmax stability check (before and after scaling)

Ans:



The screenshot shows a JupyterLab interface with a browser window at the top displaying 'localhost' and 'Q3'. The JupyterLab header indicates 'Q3 Last Checkpoint: 19 minutes ago' and 'Not Trusted'. The interface includes a menu bar (File, Edit, View, Run, Kernel, Settings, Help) and a toolbar with icons for file operations and code execution. The main area displays the results of three code cells:

A. Attention Weight Matrix

```
[8]: print("Attention Weights (Softmax):\n", attn_weights)
```

Attention Weights (Softmax):
tensor([[[0.3996, 0.2274, 0.0969, 0.1289, 0.1472],
[0.2849, 0.0494, 0.2894, 0.3574, 0.1789],
[0.0943, 0.5423, 0.1192, 0.1276, 0.1166],
[0.0767, 0.2918, 0.1180, 0.0897, 0.4238],
[0.0494, 0.6288, 0.0460, 0.0428, 0.2346]]]])

B. Output Vectors

```
[11]: print("\nOutput Vectors:\n", output)
```

Output Vectors:
tensor([[-0.5987, -0.4586, 0.1636, 0.9858],
[-0.2145, -0.6599, 0.7453, 0.2710],
[-0.5828, 0.0075, 0.4739, 1.8518],
[-0.5885, -0.5899, 0.5524, 0.8588],
[-1.0159, -0.0569, 0.4155, 1.3862]])

C. Softmax Stability Check

```
[14]: print("\nSoftmax Stability Check:")  
print("Max [raw scores] before scaling :", scores.abs().max().item())  
print("Max [scaled scores] after scaling:", scaled_scores.abs().max().item())
```

Softmax Stability Check:
Max [raw scores] before scaling : 3.958254814147949
Max [scaled scores] after scaling: 1.9791274870739746