

# Sequence-to-Sequence Models

**Tianxiang (Adam) Gao**

October 28, 2024

# Outline

1 Word Embedding

2 Sequence-to-Sequence Models

# Recap: Recurrent Neural Networks

- **Audio Waveform:** A 1D array represents the amplitude of the sound over time, e.g., 16kHz
- **One-Hot Encoding:** Each word in a **vocabulary** is a **binary** one-hot vector.
- **Challenges in Text Data:** Curse of dimensionality and long-run dependencies.
- **Language Models:** Assigns probabilities to a given sequence of words
- **Neural Language Model:** Model the probability distribution of the next word given the history:

$$\mathbb{P}(\mathbf{x}_{t+1} \mid \mathbf{x}_1, \dots, \mathbf{x}_t) = f_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_t).$$

- **RNNs:** Encode the history into a hidden state  $\mathbf{h}_t$  updated by combining with the current word  $\mathbf{x}_t$ :

$$\mathbf{h}_t = \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t + \mathbf{b}_h), \quad \text{and} \quad \hat{\mathbf{y}}_t = \text{softmax}(\mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y)$$

- **Training RNNs:** backpropagation through time
  - Forward (simplified):  $\mathbf{h}_t = \phi(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_x \mathbf{x}_t)$
  - Backward (simplified):  $d\mathbf{h}_t = \mathbf{W}_h^{\top} (\phi'_{t+1} \odot d\mathbf{h}_{t+1}) + \mathbf{W}_y^{\top} (\sigma'_t \odot d\mathbf{y}_t)$
- **Generation:** Sample the next word from the predicted probability distribution produced by RNNs.
- **RNN Types:** One-to-many, many-to-one, or many-to-many structures for different tasks.
- **Vanishing/Exploding Gradients:**  $\mathbf{h}_t = \mathcal{O}(a^t)$  and  $d\mathbf{h}_t = \mathcal{O}(b^{T-t})$ .

## Recap: Recurrent Neural Networks

- **Gated Recurrent Unit (GRU):** Gates helps maintain long-term dependencies:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h(\mathbf{r}_t \odot \mathbf{h}_{t-1}) + \mathbf{W}_x \mathbf{x}_t), \quad \text{and} \quad \mathbf{h}_t = \mathbf{z}_t \mathbf{h}_{t-1} + (1 - \mathbf{z}_t) \odot \tilde{\mathbf{h}}_t$$

- **Long Short-Term Memory (LSTM):** Use a cell state  $\mathbf{c}_t$  to maintain long-term dependencies.

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t, \quad \text{and} \quad \mathbf{h}_t = \mathbf{o}_t \odot \tanh(\mathbf{c}_t)$$

- **Bidirectional RNNs:** The concatenated hidden state:  $\mathbf{h}_t = [\vec{\mathbf{h}}_t, \overleftarrow{\mathbf{h}}_t]$
- **Deeper RNNs:** Each layer  $\ell$  computes its hidden state using the hidden state from the layer  $\ell - 1$ :

$$\mathbf{h}_t^{(\ell)} = \tanh(\mathbf{W}_h^{(\ell)} \mathbf{h}_{t-1}^{(\ell)} + \mathbf{W}_x^{(\ell)} \mathbf{h}_t^{(\ell-1)})$$

- **Drawbacks of One-Hot Representation:** orthogonality and high dimensionality
- **Word Embedding:** Words are represented as **dense vectors** in a **lower-dimensional space**.

$$\mathbf{e} = \mathbf{E}\mathbf{x}$$

- **Continuous Bag of Words (CBOW):** Predicts the target word given the context
- **Skip-Gram:** Predicts the context words given a target word.
- **Negative Sampling:** Reformulates the context-target predictions as a **binary** classification:

$$\mathcal{L}(\mathbf{E}) = - \sum_{(\mathbf{e}_c, \mathbf{e}_t)} \log \sigma(\mathbf{e}_t^\top \mathbf{e}_c) - \sum_{(\mathbf{e}_c, \tilde{\mathbf{e}}_t)} \log \sigma(-\mathbf{e}_t^\top \tilde{\mathbf{e}}_c),$$

where  $\tilde{\mathbf{e}}_t$  is **negative** target samples outside the context window.

# Outline

1 Word Embedding

2 Sequence-to-Sequence Models

# Word Analogy Task

- **Objective:** Assess the quality of word embeddings by testing how well they capture **semantic** and **syntactic** relationships between words.
- **Goal:** Given an analogy of the form: "*a is to b as c is to ???*", find the word *d* that completes the analogy correctly.
- **Examples:**
  - Semantic analogy: "*Paris is to France as Washington, D.C. is to the United States.*"
  - Syntactic analogy: "*Run is to Running as Swim is to Swimming.*"
- **Vector Arithmetic:**  $e_d \approx e_b - e_a + e_c$ , i.e., linear space
- The word *d* is chosen as the vector  $e_d$  closest to  $e_b - e_a + e_c$  based on *cosine similarity*:

$$d = \arg \max_{x \in V} S_C(e_x, e_b - e_a + e_c)$$

where  $V$  is the vocabulary, and

$$S_C(\mathbf{u}, \mathbf{v}) = \cos(\theta) = \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \cdot \|\mathbf{v}\|}$$

where  $\theta$  is the angle between vectors  $\mathbf{u}$  and  $\mathbf{v}$ .

- **Evaluation:** The predicted word *d* is evaluated by comparing it to the correct answer from the analogy dataset.

# GloVe: Global Vectors for Word Representation

- **Objective:** Create a word embedding model that captures both local context and **global** statistical information from a text corpus.

- **Co-occurrence Matrix:**

- $X_{ij}$ : Number of times word  $j$  appears in the context of word  $i$ .
- $X_i = \sum_j X_{ij}$ : Total occurrences of any word in the context of word  $i$ .
- $\mathbb{P}(w_j | w_i) = X_{ij}/X_i$ : Probability of word  $j$  occurs in the context of word  $i$ .

- **Word Comparison in Context:**

- Compare words  $e_i$  and  $e_j$  in the context  $\tilde{e}_k$  using a **probability ratio**:

$$\exp \left\{ (e_i - e_j)^\top \tilde{e}_k \right\} = \frac{\mathbb{P}(w_i | w_k)}{\mathbb{P}(w_j | w_k)} = \frac{X_{ki}}{X_{kj}}$$

- Taking the log yields:

$$e_i^\top \tilde{e}_k - e_j^\top \tilde{e}_k = \log X_{ki} - \log X_{kj} \quad \Rightarrow \quad e_i^\top \tilde{e}_k \sim \log X_{ki}$$

- **Cost Function:**

- The GloVe model learns word vectors  $e_i$  and context vectors  $\tilde{e}_k$  by minimizing:

$$J = \sum_{i,k} f(X_{ik}) \left( e_i^\top \tilde{e}_k + b_i + \tilde{b}_k - \log(X_{ik}) \right)^2$$

- $f(X_{ik})$  is a weighting function for co-occurrences, while  $b_i$  and  $\tilde{b}_k$  are bias to maintain symmetry.

## Bias in Word Embeddings

- **Problem:** Word embeddings trained on large datasets often encode societal biases, like gender stereotypes:

*"Man is to Computer Programmer as Woman is to Homemaker?"*

- **Identifying Bias Direction:** Compute the average difference between **definitional pairs**:

$$\begin{cases} \vec{e}_{\text{man}} - \vec{e}_{\text{woman}} \\ \vec{e}_{\text{he}} - \vec{e}_{\text{she}} \\ \dots \end{cases} \Rightarrow \vec{e}_{\text{bias}} = \frac{1}{N} \sum_{i=1}^N (\vec{e}_{x_i} - \vec{e}_{y_i})$$

This average can be replaced by advanced techniques like PCA.

- **Neutralization:** Project non-definitional words onto the space orthogonal to the bias direction to remove bias:

$$\vec{e} \leftarrow \vec{e} - \frac{\vec{e}^\top \vec{e}_{\text{bias}}}{\|\vec{e}_{\text{bias}}\|^2} \cdot \vec{e}_{\text{bias}}$$

- **Equalizing Pairs:** Adjust equalize pairs (like "brother" and "sister") to have equal and opposite projections along the gender direction, making them **equidistant** from the gender-neutral axis.
- **Identifying Gendered Words:** Train a classifier to distinguish between gender-specific and neutral words using a set of definitional pairs.



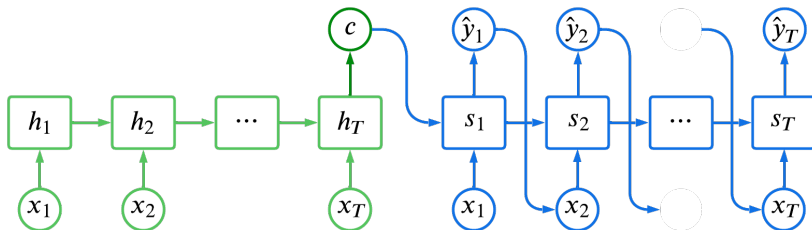
# Outline

1 Word Embedding

2 Sequence-to-Sequence Models

# Seq2Seq: Sequence-to-Sequence Models

- **Define:** Sequence-to-Sequence (Seq2Seq) models are designed to handle tasks where both input and output are sequences of variable length, e.g., machine translation or summarization.
- **Example:** “Jane visite l’Afrique en septembre.”  $\Rightarrow$  “Jane visits Africa in September.”
- **Encoder-Decoder Architecture**
  - **Encoder:** Processes and compresses the input sequence into a fixed-length **context vector**.
  - **Decoder:** Uses the context vector to generate the output sequence sequentially



Here  $c$  is the context vector summarizing the *entire* input sequence.

# Conditional Language Model

- **Language Model:** Assigns probabilities to a sequence of words  $\{\mathbf{x}_t\} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ :

$$\mathbb{P}(\{\mathbf{x}_t\}) = \mathbb{P}(\mathbf{x}_1, \dots, \mathbf{x}_T) = \prod_{t=1}^T \mathbb{P}(\mathbf{x}_t \mid \mathbf{x}_1, \dots, \mathbf{x}_{t-1})$$

where each conditional probability is modeled as:

$$\mathbb{P}(\mathbf{x}_t \mid \mathbf{x}_1, \dots, \mathbf{x}_{t-1}) = f_{\theta}(\mathbf{x}_1, \dots, \mathbf{x}_{t-1})$$

- **Conditional Language Model:** Assigns probabilities to a target sequence  $\{\mathbf{y}_t\}$  *given* an input sequence  $\{\mathbf{x}_t\}$ :

$$\mathbb{P}(\{\mathbf{y}_t\} \mid \{\mathbf{x}_t\}) = \prod_{t=1}^{T'} \mathbb{P}(\mathbf{y}_t \mid \mathbf{y}_1, \dots, \mathbf{y}_{t-1}, \mathbf{x}_1, \dots, \mathbf{x}_T)$$

where the conditional probability for each word in the target sequence is:

$$\mathbb{P}(\mathbf{y}_t \mid \mathbf{y}_1, \dots, \mathbf{y}_{t-1}, \mathbf{x}_1, \dots, \mathbf{x}_T) = g_{\phi}(\mathbf{y}_{t-1}, \mathbf{s}_t, \mathbf{c})$$

with the entire input sequence stored in the **context vector**  $\mathbf{c}$ .

# Beam Search

**Greedy Search:** Selects the most probable word at each step, which may lead to suboptimal sequences.

**Beam Search:** Tracks **multiple** high-probability sequences simultaneously to improve overall accuracy.

- **Initialization:** Start with the *seed* token and choose the top  $k$  words based on the probability distribution.
- **Expansion:** Predict the next word for each candidate, generating new sequences.
- **Pruning:** Keep the top  $k$  sequences with the highest log-probability scores, discarding the rest.

$$\mathbb{P}(\mathbf{y}_2, \mathbf{y}_1 \mid \mathbf{c}) = \mathbb{P}(\mathbf{y}_2 \mid \mathbf{c}, \mathbf{y}_1) \cdot \mathbb{P}(\mathbf{y}_1 \mid \mathbf{c})$$

- **Repeat:** Continue expanding and pruning until an *end-of-sequence* token is generated or max length is reached.

## Remark

- **Beam Width ( $k$ ):** Requires  $k$  identical decoders to update candidate sequences simultaneously.
- **Advantages:** Balances between accuracy and computation; larger  $k$  increases accuracy but demands more computational resources.

# Numerical Stability and Error Analysis

## Log-Probabilities for Stability:

- Since  $\mathbb{P}(\cdot) \in [0, 1]$ , the product of probabilities can approach zero, causing numerical instability.
- To address this, log-probabilities are used:

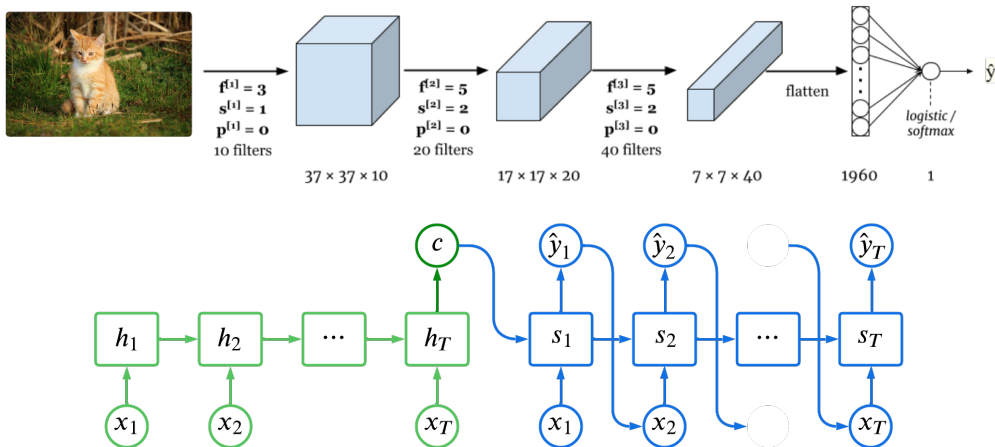
$$\begin{aligned}\mathbf{y}^* &= \underset{\mathbf{y}}{\operatorname{argmax}} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_{T'} \mid \mathbf{c}) \\ &= \underset{\mathbf{y}}{\operatorname{argmax}} \frac{1}{T'} \sum_{t=1}^{T'} \log \mathbb{P}(\mathbf{y}_t \mid \mathbf{c}, \mathbf{y}_1, \dots, \mathbf{y}_{t-1})\end{aligned}$$

This transformation helps prevent underflow and enables stable computation of probabilities.

## Error Analysis in Beam Search:

- Let  $\mathbf{y}^*$  be the optimal sequence and  $\hat{\mathbf{y}}$  the model's predicted sequence.
- **If**  $\mathbb{P}(\mathbf{y}^* \mid \mathbf{c}) > \mathbb{P}(\hat{\mathbf{y}} \mid \mathbf{c})$ : Increasing beam width can improve accuracy by exploring more potential sequences.
- **If**  $\mathbb{P}(\mathbf{y}^* \mid \mathbf{c}) \leq \mathbb{P}(\hat{\mathbf{y}} \mid \mathbf{c})$ : Increasing beam width won't help, as errors stem from model limitations.

## Image Captioning



*"A small orange kitten sits attentively on green grass, surrounded by natural, dried foliage in the background, giving a calm and serene outdoor setting."*

# BLEU Score: Bilingual Evaluation Understudy

**BLEU Score:** Bilingual Evaluation Understudy (BLEU) evaluates a machine-generated **candidate translation**  $\hat{y}$  by comparing it to a *list* of reference translations  $\{y_1, \dots, y_M\}$ .

- **Candidate Translation:** "the cat the cat sat"
- **Reference Translation 1:** "the cat sat on the mat"
- **Reference Translation 2:** "the cat sat by the mat"

**Precision:** Measures how many n-grams in the candidate match the reference translations.

- **Candidate Bigrams:** {"the cat", "cat the", "the cat", "cat sat"}
- **Reference 1 Bigrams:** {"the cat", "cat sat", "sat on", "on the", "the mat"}
- **Reference 2 Bigrams:** {"the cat", "cat sat", "sat by", "by the", "the mat"}

$$\text{Precision} = \frac{\text{Number of matching n-grams}}{\text{Total n-grams in candidate}} = \frac{3}{4} \implies \text{Inflate the score!}$$

## Modified Precision

**Modified Precision:** Limits the count of an n-gram to the maximum it appears in any reference.

- **Count in Candidate:**

- "the cat" appears **2 times** in the candidate.
- "cat the" appears **1 time** in the candidate.
- "cat sat" appears **1 time** in the candidate.

- **Clipped Count:**

- "the cat" appears at most **1 time** in any reference.
- "cat sat" appears at most **1 time** in any reference.

$$\text{Modified Precision} = \frac{\text{Clipped number of matching n-grams}}{\text{Total n-grams in candidate}} = \frac{2}{4}$$

**BLEU Score Formula:** Uses modified precision  $p_n$  for n-grams:

$$\text{BLEU} = BP \cdot \exp \left( \frac{1}{N} \sum_{n=1}^N \log p_n \right)$$

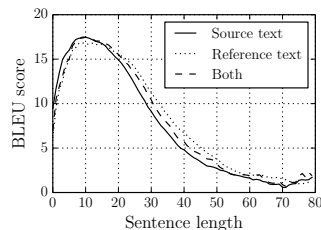
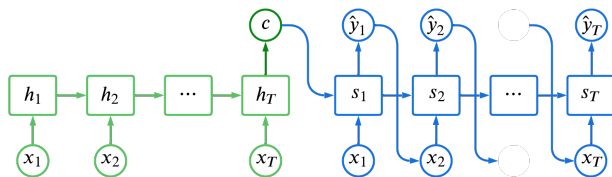
- $N$ : Maximum n-gram length considered (typically  $N = 4$ ).
- **Brevity Penalty (BP):**

$$BP = \begin{cases} \exp \left( 1 - \frac{r}{c} \right), & \text{if } c < r \\ 1, & \text{if } c \geq r \end{cases}$$

where  $c$  is the length of the candidate translation and  $r$  is the length of the reference translation.



# Limitations of RNN Encoder-Decoder Framework



- **Encoder:** Process and compress the input sequence  $\{x_1, \dots, x_T\}$  into a context vector  $c$ .

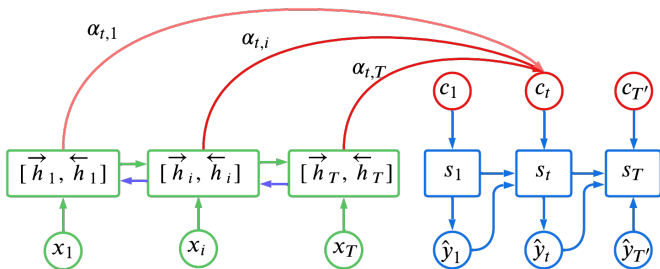
$$h_t = f_{\theta}(h_{t-1}, x_t), \quad c = h_T.$$

- **Decoder:** Uses the context vector  $c$  to generate the output sequence sequentially

$$\mathbb{P}(y_t | x, y_1, \dots, y_{t-1}) = \mathbb{P}(y_t | s_t), \quad \text{where } s_t = g_{\phi}(s_{t-1}, y_{t-1}, c)$$

**Note:** Encoding all information into a **single** vector  $c$  may cause information loss for longer sequences.

## Distinct Context Vector in Attention Mechanism



- **Distinct Context Vector for Each Target Word:** Each target word  $y_t$  has a unique context vector  $c_t$ , allowing the model to focus on relevant input parts.

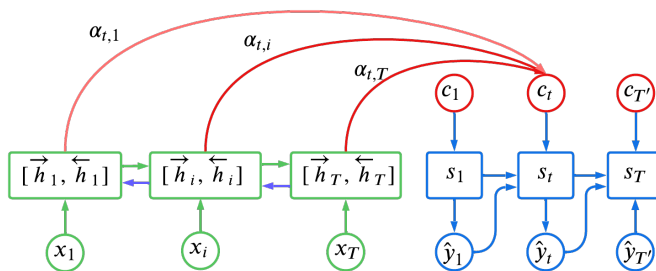
$$\mathbb{P}(y_t \mid x, y_1, \dots, y_{t-1}) = \mathbb{P}(y_t \mid s_t), \quad \text{where } s_t = g_\phi(s_{t-1}, y_{t-1}, \mathbf{c}_t)$$

- **Context Vector Computation:** The context vector  $c_t$  is computed as a weighted sum of encoder hidden states  $h_i$ , tailored to the current decoding step.

$$\mathbf{c}_t = \sum_{i=1}^T \alpha_{t,i} \mathbf{h}_i$$

where **attention weights**  $\alpha_{t,i}$  indicate the relevance of each hidden state  $h_i$  for generating  $y_t$ .

# Attention Mechanism in Seq2Seq



- **Attention Weights:** Computed from an **alignment score**  $e_{t,i}$  between the decoder's previous hidden state  $s_{t-1}$  and each encoder hidden state  $h_i$ .

$$\alpha_{t,i} = \text{softmax}(e_{t,i}), \quad \text{where} \quad e_{t,i} = a(s_{t-1}, h_i)$$

- **Alignment Model:** The alignment function  $a$  is a feed-forward neural network, trained jointly with Seq2Seq to optimize attention.
- **Bidirectional Encoder:** The encoder uses a bidirectional RNN to capture both past and future context, enhancing comprehension of each input word's meaning.