

Instruction-Tuning

Estimated Time: 45 minutes

Learning objectives

After completing this LAB, you will be able to:

- Explain instruction marking and why it is used in training transformer models
- Explain how instruction marking improves model performance by focusing on response tokens
- Identify the role of loss function minimization in instruction marking
- Explain how instruction marking is applied in a context or other real-world scenarios
- Analyze the impact of cross-entropy loss in training models with instruction marking

Introduction

As transformer-based models evolve, fine-tuning the models such as instruction-tuning have become critical for improved response quality in AI systems. However, traditional loss functions often treat all tokens in a sequence equally, which can lead to inefficiencies in learning. Instruction marking addresses this by selectively focusing loss calculation on response tokens while ignoring instruction tokens. This ensures that models prioritize generating accurate and contextually relevant responses.

Instruction marking

Instruction marking is a sophisticated technique employed during instruction-tuning to enhance the learning process of transformer models. By focusing the loss calculation on specific parts of the output sequence – such as response tokens – while ignoring instruction tokens, instruction marking ensures that the model prioritizes generating accurate and contextually relevant responses.

Why use instruction marking?

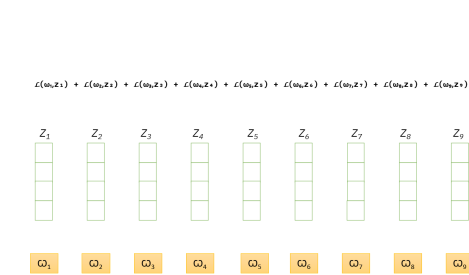
- **Focus on response:** The model learns to generate precise and meaningful responses.
- **Efficient learning:** Reduces computational overhead by excluding irrelevant tokens from loss calculations.
- **Promotes overfitting:** This minimizes the risk of the model overfitting to instruction tokens, which are often repetitive or less informative.

Instruction marking and loss function representation

Loss function representation

The total loss function for instruction-tuning is expressed as the sum of individual loss functions for each token pair:

$$L(y_1, x_1) + L(y_2, x_2) + L(y_3, x_3) + \dots + L(y_n, x_n)$$



The image represents elements involved in instruction marking, showing a structured relationship between

• **Input:** Instruction tokens (x_1, x_2, \dots, x_n) and response tokens (y_1, y_2, \dots, y_n).
• **Output:** A sequence of loss functions ($L(x_i, y_i)$) calculated for each token pair, indicating an optimization process for learning objectives tied to token elements.

$$L(x_1, y_1) + L(x_2, y_2) + L(x_3, y_3) + L(x_4, y_4) + L(x_5, y_5) + L(x_6, y_6) + L(x_7, y_7) + L(x_8, y_8) + L(x_9, y_9)$$

Where:

- L : Loss function
- n : (integer) Batch size/instruction length
- i : (integer) Token index/instruction

Visual components

Vertical bars ($i=1$ to n): Marked positions.

Yellow boxes ($n=1$ to n): Input context.

Combined pairs: Processed token pairs for calculation.

How does the instruction marking work?

To understand how instruction marking works in a real-time scenario, let's break it down step by step using a practical example.

Real-time scenario: Chatbot for answering questions

A chatbot is trained to answer user queries using instruction-tuning. Instruction marking ensures that the model prioritizes response generation over memorizing context.

Example: Input sequence

Instruction: Which is the largest ocean?
Response: The Pacific Ocean.

Here, the input consists of two parts:

- **Instruction:** The user's question or command (Which is the largest ocean?).
- **Response:** The expected answer generated by the model (The Pacific Ocean).

Marking process

In instruction-tuning, the goal is to ensure that the model learns only from the response tokens while ignoring the instruction tokens. Here's how it works:

- **Instruction tokens are masked (loss ignored):**
 - ➔ Tokens from the instruction (e.g., Which, is, the, largest, ocean,) are ignored during loss calculation.
 - ➔ Prevents overfitting to instruction formats.
- **Response tokens remain unmasked (loss calculated):**
 - ➔ Tokens from the answer (e.g., The, Pacific, Ocean, is) are included in loss calculation.
 - ➔ Trains the model to generate accurate, context-aware responses.

Processing steps

Let's walk through the steps involved in processing the input sequence with instruction marking:

- **Tokenization:**

The input sequence is split into tokens (words or words) using a tokenizer.

Example:

Instruction: "Which, is, the, largest, ocean,?"
Response: "The, Pacific, Ocean,."

- **Encoding:**

Each token is converted into a numerical representation (embedding) that the model can process.

Example:

"Which" → [0.23, 0.45, ...]
"Pacific" → [0.45, 0.89, ...]

- **Position embedding:**

Positional information is added to the tokens to preserve the order of the sequence.

Example:

"Which" (Position 1) → [0.23, 0.45, ...] + [0.1, 0.2, ...]
"Pacific" (Position 7) → [0.45, 0.89, ...] + [0.8, 0.9, ...]

- **Context preservation:**
 - ➔ The model processes the tokens through its layers (e.g., transformer layers) to capture contextual relationships.

Example:

The model understands that "largest" is related to "ocean" and "Pacific".

- **Token prediction:**
 - ➔ The model predicts the next token in the sequence.

Example:

After processing the instruction, the model predicts "The Pacific Ocean" as the next token.

- **Loss calculation:**

The model calculates the loss for each predicted token against the valid tokens.

Example:

For the token "Pacific", predicted with 0.7 probability: $L(w_i, x_i) = -\log(0.7) \approx -(0.357) \approx 0.357$

For an incorrect prediction (e.g., "Atlantic") with 0.2 probability: $L(w_i, x_i) = -\log(0.2) \approx -(1.609) \approx 1.609$

- **Masking application:**
 - ➔ The model applies instruction marking to focus loss calculation on response tokens only.

Example:

Masked tokens: "Which, is, the, largest, ocean,?"
Unmasked tokens: "The, Pacific, Ocean,."

Cross-entropy loss in instruction marking

Cross-entropy loss is commonly used in training models with instruction marking.

It measures how well the model predicts the correct output by comparing the predicted probability distribution with the proper labels. The marked tokens (response parts of the input) are the focus in the loss calculation, preventing them from affecting learning.

Cross-entropy loss (Mathematical Formulation)

The cross-entropy loss measures how well the model's predicted probabilities match the true distribution. For a single token, it is defined as:

$$\text{Cross-Entropy Loss} = -\sum_{i=1}^V y_i \log(p_i)$$

Where:

- y_i : True probability distribution (one-hot encoded, 1 for the correct token, 0 for others).
- p_i : Predicted probability for token i .
- V : Vocabulary size.

Illustrative example: Question answering

We have a sequence of ground truth tokens (correct words in a sentence) and predicted tokens (what the model generated).

Example values:

"Which is the largest ocean? The Pacific Ocean."

Let's break down the concepts with the same example:

- **Input:** The input and the expected answer are sequences of tokens (words or sub-words). For instance:
 - ➔ Input (Instruction Tokens): "Which, is, the, largest, ocean,?"
 - ➔ Expected Answer (Answer Tokens): "The, Pacific, Ocean,."
- **Model prediction:** The model generates predictions (\hat{y}_i) for each token in the expected answer. These predictions are probability distributions over the vocabulary.
- **Ground truth:** The ground truth (y_i) is a one-hot vector indicating the correct token at each position in the answer.

This document consists of input, an instruction, and an answer.

Index	Ground Truth Token (x)	Model Prediction (y)
x ₁	Where	y ₁
x ₂	is	y ₂
x ₃	the	y ₃
x ₄	largest	y ₄
x ₅	ocean?	y ₅
x ₆	The	y ₆
x ₇	Pacific	y ₇
x ₈	Ocean	y ₈
x ₉		y ₉

- The first five tokens (x₁ to x₅) are **instruction tokens** (question part).
- The last four tokens (x₆ to x₉) are **answer tokens** (response part).
- The model generates predictions (y₁, y₂, ..., y₉) corresponding to each token.

Example calculation:

We are dealing with a sequence prediction problem, where a model predicts words or tokens, and we compute how good or bad its predictions are using a loss function.

Step 1: Understanding the input and prediction

Input:

- "Where is the largest ocean?" (x₁ to x₅)

Possible answer probabilities (as given by the model):

"Pacific Ocean" → 0.7 (Correct Answer)

"Atlantic Ocean" → 0.2

"Indian Ocean" → 0.1

Since the correct answer is "Pacific Ocean", the model should have assigned a probability as close to 1.0. However, it gave 0.7.

Step 2: Understanding the function

We use the negative log-likelihood (NLL) loss function, defined as:

$$Loss_{NLL} = -\log(P_{\text{pred}}(x))$$

This function measures how "wrong" the model's predicted probability is for the correct answer. The logarithm function is used because:

- It penalizes lower probabilities more.
- It ensures that probabilities close to 1 have a low loss (indicating a good prediction).
- It is a standard choice in classification problems (often associated with cross-entropy loss).

Step 3: Applying the formula

Since the model predicted P("Pacific Ocean") = 0.7, we substitute:

$$Loss_{NLL} = -\log(0.7)$$

Using the natural logarithm:

$$\log(0.7) \approx -0.357$$

So,

$$Loss_{NLL} \approx -(-0.357) = 0.357$$

Interpretation

- A lower cross-entropy loss means the model is confident about the correct answer.
- If the "Pacific Ocean" probability had been higher, the loss would be lower.
- If the "Pacific Ocean" probability were lower (e.g., 0.5 instead of 0.7), the loss would be higher, indicating the model is further from correct.

Summary table

Prediction Probability	Log Calculation	Final Loss
0.7 (Correct)	$-\log(0.7)$	0.357
0.2 (Wrong)	$-\log(0.2) \approx -(-1.609)$	1.609
0.1 (Wrong)	$-\log(0.1) \approx -(-2.303)$	2.303
0.5 (Very Bad)	$-\log(0.5) \approx -(-0.693)$	0.693

This table shows that higher certainty in correct predictions results in a lower loss, while lower certainty leads to higher loss.

Conclusion

Integrating these advanced transformer-based models into existing loss computation pipelines offers numerous benefits, including improved accuracy and contextual relevance. This thorough analysis reveals the underlying mechanisms, optimizes computational efficiency, and enhances your understanding of AI-driven applications such as chatbots and QA systems. As AI evolves, continuous learning will remain critical in refining models for real-world usability.

Key Takeaways

- Profound Learning** - Loss computation isn't just about finding errors; it's about understanding model accuracy and adaptability.
- Better Generalization** - Prevents overfitting to specific data patterns, making models more adaptable.
- Optimized Computation** - Balances training overhead while maintaining superior quality.
- Context Awareness** - Maximizes contextual understanding while refining response generation.
- Real-World Impact** - Improves AI application like chatbots and QA systems.

Author(s)

- Sophia Chen

Other Contributor(s)

- Michael Smith
- Liam O'Reilly



Skills Network