

Title: Image Captioning using Two-Stage Deep Neural Network: A Performance Analysis

Abstract: This report presents the design and evaluation of a two-stage deep neural network pipeline for image captioning using the Flickr8K dataset. The pipeline comprises a Convolutional Neural Network (CNN) Encoder for image processing, an LSTM Decoder for caption generation, and a Feedforward Neural Network (FFNN) for result refinement. The performance of the model is assessed using the BLEU score, achieving an accuracy of 0.86 as the best similarity score.

1. Introduction: Image captioning is a complex task that involves understanding visual content and generating coherent natural language descriptions. This report presents a comprehensive pipeline designed to tackle this task using a two-stage deep neural network approach.

2. Methodology: The pipeline consists of three main components:

2.1 CNN Encoder: A Convolutional Neural Network (CNN) is employed to encode the image features. The encoder extracts high-level visual representations from the input images, which are then used as a basis for generating captions. The CNN's architecture is tailored for efficient feature extraction while maintaining spatial information.

2.2 LSTM Decoder: A Long Short-Term Memory (LSTM) Decoder is utilized to generate captions based on the encoded image features. The LSTM network is trained to predict each word in the caption sequentially, considering the context of the previous words. This enables the model to capture the semantic relationships within the caption.

2.3 FFNN Refinement: A Feedforward Neural Network (FFNN) is integrated into the pipeline to refine the generated captions. This network fine-tunes the output of the LSTM decoder, enhancing the overall quality and coherence of the captions. The FFNN's architecture is optimized to capture nuances in language and improve the fluency of the generated text.

3. Dataset: The pipeline is evaluated using the Flickr8K dataset, which contains a diverse collection of images paired with human-generated captions. The dataset provides a substantial variety of scenes, objects, and contexts, making it suitable for assessing the model's performance.

4. Performance Evaluation: The performance of the pipeline is quantified using the BLEU (Bilingual Evaluation Understudy) score. BLEU measures the similarity between the generated captions and the reference captions in the dataset. The higher the BLEU score, the better the generated captions match the reference captions.

5. Results: Upon evaluation, the pipeline achieves a notable accuracy of 0.86 as the best similarity score using the BLEU metric. This indicates a strong alignment between the generated captions and the reference captions, showcasing the efficacy of the two-stage deep neural network approach.

6. Conclusion: In this report, we have demonstrated the successful implementation of a two-stage deep neural network pipeline for image captioning. The pipeline utilizes a CNN Encoder, an LSTM Decoder, and an FFNN Refinement stage to generate accurate and coherent captions for images. The achieved accuracy of 0.86, as measured by the BLEU score, underscores the effectiveness of the approach in producing high-quality captions.

7. Future Work: Future work could involve exploring alternative neural network architectures, experimenting with different pre-processing techniques, and investigating additional evaluation metrics to gain a more comprehensive understanding of the model's performance.

Acknowledgments: We acknowledge the support of the Flickr8K dataset creators for providing the dataset that facilitated this research.