

ASSIGNMENT STATUS REPORT

L02: Exploring Deep Learning Tools - A No-Code Introduction to TensorFlow and Keras

This lab was a great way for me to get hands-on experience with deep learning concepts using the pre-trained VGG16 model, all without needing to write any code myself. The introduction set the stage well, explaining how deep learning models function and their application in image classification. We used the VGG16 model, a well-known architecture trained on the ImageNet dataset, to classify images based on the features it has learned. This really helped solidify my understanding of pre-trained models and what they can do.

First, I loaded the VGG16 model and spent some time looking at its architecture. By interacting with the pre-loaded code cell, I could see the model's complex structure, with its multiple convolutional and fully connected layers. It reinforced the idea that deep learning is heavily about feature extraction, with each layer contributing to the image classification process. Understanding this architecture made me appreciate the complexity and efficiency of these models in recognizing objects even more.

Next, I focused on data loading and preprocessing. The read-only code cell clearly showed how crucial it is to prepare images before feeding them into the model. This preprocessing involved resizing images to match the VGG16 model's required input dimensions and normalizing pixel values for consistency. It was clear that proper preprocessing is essential in deep learning; otherwise, you risk inaccurate predictions or poor model performance.

The most interactive part of the lab was making predictions with the VGG16 model. I uploaded three images: a cat, a dog, and a deer, and observed how the model classified each. The results were quite interesting!

- The cat image was misclassified; the model predicted 'hamper' with 16.16% probability, followed by 'carton' (15.81%) and 'tub' (5.90%). This showed the model's limitations when dealing with objects not prominently featured in its training data.
- The dog image was classified more accurately. The top prediction was 'Labrador Retriever' (42.55%), then 'Golden Retriever' (27.97%), and 'Chesapeake Bay Retriever'

(6.79%). This demonstrated VGG16's ability to be quite specific in classifying common dog breeds.

- The deer image was classified as 'gazelle' (44.69%), with 'impala' (16.38%) and 'Mexican hairless' (9.49%) as secondary predictions. This was reasonably close, as gazelles and impalas do look similar to deer.

This experiment also highlighted how sensitive deep learning models can be to minor variations in test data. It underscored the importance of data augmentation during training to improve model robustness and generalization.

The final part of the lab was a reflection on these insights. It demonstrated that while pre-trained models like VGG16 perform decently in image classification, they aren't perfect. Misclassifications can occur, especially with objects outside the training distribution or when input images are distorted. Another key takeaway was the critical role of preprocessing in ensuring prediction accuracy.

Even though the cat image I used was perfectly clear, I was surprised that the model got it completely wrong, while the predictions for the dog and deer were quite good.

Overall, I found this lab very informative. From loading the model to making predictions and analyzing the results, it provided a solid introduction to deep learning. The interactive elements really helped reinforce theoretical concepts with practical, hands-on experience. In the future, I'd be interested in exploring other pre-trained models and experimenting with fine-tuning to see if I can improve classification accuracy on specific datasets. This experience definitely deepened my understanding of deep learning models and their real-world applications.

L03: AWS MLU Lab Reflection (Module 01)

Lab 01: Getting Started with PyTorch

- **Summary of Key Learnings:** In this lab, I dove into the foundational concepts of PyTorch, with a focus on tensor manipulation and building a very basic neural network. I learned to create tensors of different shapes and data types, how to index and slice them to access specific parts, and how to perform both element-wise and matrix operations like addition, multiplication, and transpose. A key aspect was converting tensors to NumPy arrays and standard Python scalars, which is important for integrating PyTorch with the broader Python ecosystem. Following that, I generated a synthetic regression dataset, defined a simple two-layer fully connected network (using `nn.Linear`, ReLU activation, then another `nn.Linear`), and implemented the full training loop: the forward pass, loss computation (using `nn.MSELoss`), the backward pass (`loss.backward()`), and updating parameters using Stochastic Gradient Descent (SGD).
- **Insights & Understanding:**

- **Dynamic Computation Graphs:** I really came to appreciate PyTorch's eager execution model, where the computation graph is built dynamically as operations are performed. Being able to print intermediate activations and gradients while stepping through the training process really demystified backpropagation for me.
- **Autograd Mechanics:** Looking at the `.grad` attributes of tensors after the backward pass clearly showed how gradients flow from the loss all the way back to each weight in the network.
- **Data Conversion:** The process of converting between CPU and GPU tensors, NumPy arrays, and Python scalars highlighted the importance of writing device-aware code. It also made me realize how crucial it is to minimize unintended data transfers, which can lead to subtle bugs and performance issues.
- **Challenges Encountered:**
 - **Shape Mismatches:** Initially, my network definition had mismatched input/output dimensions between layers, which caused runtime errors. I fixed this by carefully printing the shapes of tensors before each layer to ensure they were consistent.
 - **Learning Rate Sensitivity:** I found that a learning rate that was too large (e.g., 0.1) caused the loss to diverge. Reducing it to 0.01 stabilized the training. This really drove home how important hyperparameter tuning is, even for simple models.
- **Application & Relevance:** Understanding low-level tensor operations is directly applicable to building custom layers, like attention mechanisms, where I'd need to manipulate multi-dimensional arrays. Debugging training loops at this fundamental level prepares me for more complex projects, such as semantic segmentation on Mars imagery, where forward and backward passes might involve custom operations.
- **Code and Experimentation:**
 - I experimented with different learning rates (0.1, 0.01, 0.001) and found that 0.01 converged the fastest without causing oscillations.
 - I also tried adding an extra hidden layer, which doubled the model's capacity. This reduced the training loss but led to some mild overfitting on the synthetic dataset I was using.

Lab 02: How Neural Networks Learn

- **Summary of Key Learnings:** This lab built upon the previous one by having me construct a multi-layer perceptron (MLP) with two hidden layers and integrate dropout for regularization. I defined this network using `nn.Sequential`, trained it on a classification task (similar to MNIST using synthetic digits), and meticulously tracked both training and validation accuracy across epochs. A key takeaway was understanding *why* overfitting occurs—models essentially memorize the training data—and how dropout helps by randomly deactivating neurons during training, which forces the network to learn more robust features.
- **Insights & Understanding:**
 - **Bias-Variance Tradeoff:** Watching the validation accuracy plateau or even start to decline while the training accuracy kept improving was a clear illustration of how overly complex models tend to overfit.

- **Dropout Effectiveness:** Introducing dropout with a probability $p=0.2$ in the hidden layers resulted in a noticeable improvement in validation accuracy (for instance, from around 85% to 89% at epoch 20). This confirmed dropout's effectiveness as a regularization technique.
- **Model Capacity vs. Data:** I observed that smaller networks tended to underfit (showing low training and validation accuracy), whereas larger networks would overfit if I didn't use regularizers.
- **Challenges Encountered:**
 - **Implementing Dropout Correctly:** My first attempt involved applying dropout during the evaluation phase as well, which unfortunately degraded the test performance. I learned that dropout layers need to be wrapped or managed such that they are only active during `model.train()` mode and automatically disabled during `model.eval()`.
 - **Monitoring Metrics:** The provided lab code only printed the loss, so I had to write custom loops to compute per-epoch accuracy. This reinforced the good practice of looking at both loss and accuracy for a holistic evaluation of the model's performance.
- **Application & Relevance:** Dropout and the design of multi-layer networks are crucial for real-world tasks such as object detection or semantic segmentation, especially when dealing with limited labeled data where overfitting is a significant concern. The principles I learned here will directly inform how I design segmentation models for projects like Mars rovers, ensuring they can generalize across diverse Martian terrains.
- **Code and Experimentation:**
 - I tested different dropout probabilities (0.1, 0.2, 0.5) and found that $p=0.2$ gave the optimal performance, balancing regularization effectively without causing underfitting.
 - I also tried increasing the width of the hidden layers from 128 to 256 neurons. When combined with dropout, this sped up training but also required implementing early stopping to prevent overfitting.

Lab 03: Building an End-to-End Neural Network Solution (First Example of Neural Networks)

- **Summary of Key Learnings:** In Lab 03, I tackled an end-to-end text data processing task. The steps included importing a CSV file containing text and multi-class labels, tokenizing the text using a vocabulary built from the dataset itself, converting these tokens into padded index sequences, and then defining an embedding-LSTM classifier. I trained this network, validated its performance on a held-out dataset, and experimented with various hyperparameters like embedding dimension, hidden layer size, and the number of training epochs. This lab really brought together data preprocessing, model building, training, and validation into one cohesive pipeline.
- **Insights & Understanding:**
 - **Text Preprocessing:** Seeing the impact of different tokenization strategies (word-level vs. subword-level) and the necessity of padding sequences really instilled an appreciation for the complexities of handling variable-length text data.
 - **Embedding Layers:** Visualizing the learned embeddings (using PCA) was fascinating. I could see that semantically similar words naturally clustered together,

demonstrating how embedding layers capture meaning.

- **Sequential Models:** Training an LSTM on text data was a distinct experience compared to working with feedforward networks. Managing sequence length, hidden states, and batch padding introduced new layers of complexity.
- **Challenges Encountered:**
 - **Handling OOV Tokens:** Rare words not present in the vocabulary caused index errors. I solved this by adding a special <UNK> (unknown) token to the vocabulary and mapping all unseen words to this token.
 - **Sequence Length Tradeoffs:** Using longer maximum sequence lengths marginally improved accuracy (by about 2%) but at the cost of increased memory usage and training time. I eventually settled on a sequence length of 100 as a good balance for efficiency.
- **Application & Relevance:** This end-to-end workflow closely mirrors real AI projects, which typically involve data ingestion, cleaning, model building, training, and evaluation. Text models are directly relevant to tasks I might encounter, like sentiment analysis of NASA mission logs or generating summaries from rover-generated text. The modular pipeline I built can also be adapted for other types of data, like images, audio, or even multimodal data.
- **Code and Experimentation:**
 - I changed the embedding dimension from 50 to 100 and observed about a 1.5% accuracy gain, though it doubled the training time.
 - I added dropout with $p=0.3$ to the LSTM outputs, which helped reduce overfitting and raised the validation accuracy from 78% to 82%.
 - I also tried swapping the LSTM for a 1D-convolutional classifier. While the training was faster, the final accuracy on text classification was about 5% lower.

L04: Lab ITAI 2376 - Convolutional Neural Networks (CNNs) with MNIST

This lab on Convolutional Neural Networks (CNNs) using the MNIST dataset was a fantastic way to get practical experience with image classification in deep learning. It really helped solidify new concepts for me and also built upon what I already knew about neural networks.

- **Learning Insights:**
 - **New CNN Concepts:** I got a practical introduction to the core components of a CNN. I now have a much clearer understanding of how convolutional layers work by sliding filters across images to extract features. I also learned how max pooling functions to aggregate features and reduce dimensionality (downsampling). The idea of a "receptive field"—how each neuron in a convolutional layer "sees" only a part of the input image—became much more concrete. It was fascinating to grasp how these

layers work together to build hierarchical representations of image data, starting from basic edges and textures in the initial layers and moving to more complex patterns and shapes in deeper layers. The lab also stressed the importance of data preparation steps like normalization and one-hot encoding.

- **Connection to Previous Knowledge:** My prior understanding of neural networks—concepts like gradient descent, activation functions, and backpropagation—provided a solid foundation for this lab. The core idea of training a model by adjusting weights to minimize a loss function was already familiar. This lab showed me how CNNs extend these fundamental principles and adapt them specifically for image data. My familiarity with dense layers in traditional neural networks made it easier to understand the purpose of the flatten layer in transitioning from convolutional layers to fully connected layers.
- **Surprises with CNNs and MNIST:** I was initially quite surprised by how well a relatively simple CNN architecture performed on the MNIST dataset. It was really satisfying to see the model's accuracy improve with each training epoch. I also realized that CNNs are significantly more computationally efficient for image tasks compared to fully connected networks. The automatic learning of features, rather than manual feature engineering, was also remarkable.
- **Challenges and Growth:**
 - **Specific Challenges:** One of my first hurdles was understanding the correct input shape required for the convolutional layers. I was also initially unsure about selecting the most appropriate optimizer and loss function for this specific multi-class classification task. Tuning hyperparameters like the batch size and the number of epochs was another challenge; I wasn't sure how to systematically adjust these or what values to start with.
 - **Overcoming Challenges:** To understand the input shape requirements for Conv2D layers, I consulted the Keras documentation and online tutorials. I learned that the input should be a 4D tensor (number of images, image height, image width, number of channels). For the optimizer and loss function, after researching common practices for multi-class classification, I found that "adam" and "categorical_crossentropy" are often good starting points. For hyperparameter tuning, I experimented with different values, starting with common defaults and then observing their impact on the model's performance. I realized how crucial it is to monitor both training and validation metrics to identify issues like overfitting or underfitting.
 - **Helpful Resources:** The Keras documentation was invaluable for understanding the details of each layer and function. Blog posts and online tutorials with step-by-step examples were also very helpful. Discussing concepts with peers helped clarify some of my doubts. Visualizing the convolutional and pooling operations, perhaps by looking at feature maps or using diagrams, helped me better understand how these layers process image input.
- **Personal Development:**
 - **Changed Understanding of Deep Learning:** This lab transformed my understanding of deep learning from a purely theoretical concept into a practical skill. Building and training a CNN from scratch has given me a much greater appreciation for the power and versatility of deep learning models. I now have a better grasp of the design

decisions involved in creating a CNN and feel more confident tackling image classification tasks.

- **Aspects to Explore Further:** I'm eager to learn more about advanced CNN architectures like ResNet and Inception, which address the challenges of training very deep networks. I also want to explore how CNNs can be applied to other computer vision tasks, such as object detection, image segmentation, and image captioning. Furthermore, I'm curious about transfer learning and how fine-tuning pre-trained models can accelerate development for new tasks.
- **New Perspectives (if already familiar):** Although I had some prior theoretical knowledge of CNNs, this hands-on application was incredibly beneficial. It allowed me to experiment with different architectures and parameters and directly observe their impact on model performance. This experience reinforced my understanding and made me realize the importance of real-world experimentation in deep learning. It also improved my appreciation for the challenges of hyperparameter tuning and the significance of thorough model evaluation.

The model I built achieved a test accuracy of 0.9912 (or 99.12%). I was very happy with this result, as it's quite high for the MNIST dataset and indicates the model learned effectively.

L05: AWS Machine Learning University Module 2 Lab Exploration (Labs 1-4)

Lab 1: Processing Text

- **Summary of Key Learnings:**
 - **Word Clouds & Frequency Analysis:** I learned to generate word clouds using the wordcloud library. After cleaning and tokenizing a corpus of about 10,000 customer reviews, I experimented with different sets of stopwords (NLTK's default plus a custom list). Visualizing the top 100 terms helped me quickly identify domain-specific jargon like "rover," "dust storm," and "sol".
 - **Stemming vs. Lemmatization:** I compared the outputs from NLTK's PorterStemmer and WordNetLemmatizer. While stemming correctly reduced "processing" to "process," it also incorrectly changed "analysis" to "analysi," which could harm interpretability. Lemmatization, on the other hand, produced valid dictionary forms and actually improved the accuracy of a downstream logistic regression model by about 3 percentage points in a small test.
 - **Part-of-Speech (POS) Tagging:** Using NLTK's pos_tag, I filtered tokens based on their part of speech (e.g., keeping only nouns and adjectives). This simple heuristic reduced the vocabulary size by 40% and sped up the vectorization process.
 - **Named Entity Recognition (NER):** I leveraged spaCy's en_core_web_sm pipeline to extract PERSON, ORG, and GPE (geopolitical entities) from a dataset of news headlines. I found that including the count of GPEs as a feature improved a headline sentiment classifier's F1-score by about 2 points.
- **Insights & Understanding:**

- **The Power of Preprocessing:** Initially, I tried feeding raw text directly into a classifier and saw less than 60% accuracy. However, after applying stemming, POS filtering, and stopword augmentation, the performance jumped to around 75%. This really drove home the point that even simple NLP preprocessing pipelines can lead to significant gains before any complex modeling is even attempted.
- **Trade-offs in Tokenization:** Experimenting with character-level bigrams versus word-level tokens showed that character bigrams captured morphological variants (like "drive" vs. "driving") well, but at the cost of a much larger feature count. In my final pipeline, I settled on word unigrams and bigrams with a minimum document frequency (`min_df`) of 5, which balanced coverage and efficiency.
- **Challenges Encountered:**
 - **Environment & Dependencies:** Installing spaCy models (e.g., `python -m spacy download en_core_web_sm`) sometimes failed due to my institution's proxy server. I managed to resolve this by setting the `HTTPS_PROXY` and `HTTP_PROXY` environment variables.
 - **Memory Constraints:** Loading a 500MB corpus into memory for word cloud generation caused my laptop to struggle with swapping. I rewrote the data loader to stream one file at a time and update token counts incrementally, which eliminated these memory spikes.
 - **POS Tagging Speed:** NLTK's POS tagger was quite slow when processing over 100,000 tokens. I experimented with spaCy's `nlp.pipe()` batch mode, which turned out to be about 5 times faster.
- **Application & Relevance:**
 - **Real-World Feature Engineering:** In industrial settings, I've seen that a robust text preprocessing pipeline often matters more than the model architecture itself. This lab reinforced the importance of creating reusable and scalable text pipelines, for example, using spaCy's Matcher patterns to generalize entity extraction.
 - **Future Projects:** For my capstone project, which involves analyzing rover telemetry logs, I plan to incorporate NER to automatically detect location mentions and anomaly keywords.
- **Code & Experimentation:**
 - **Token-Frequency Sampling:** I wrote a small script to sample 10% of the corpus for rapid prototyping, and then validated the pipeline on the full dataset once it stabilized.
 - **Hyperparameter Sweeps:** I performed sweeps for bigram weights (e.g., `ngram_range=(1,2)` with different importance levels) and observed diminishing returns past a certain point (0.7), which informed my final parameter choices.

Lab 2: Using the Bag-of-Words Method

- **Summary of Key Learnings:**
 - **CountVectorizer vs. TfidfVectorizer:** I built document-term matrices using scikit-learn's CountVectorizer (with `max_features=5000`, `min_df=10`) and TfidfVectorizer. On a spam dataset of 20,000 documents, raw counts yielded about 82% accuracy with logistic regression, while TF-IDF pushed this up to around 88%.
 - **Binary vs. Frequency vs. TF-IDF Representations:** I compared three different modes: binary presence/absence of terms, raw term counts, and TF-IDF scores. Binary

features led to faster training (about 30% less time) but resulted in slightly lower accuracy (around 85%) compared to TF-IDF's ~88%.

- **N-gram Configurations:** Including bigrams (`ngram_range=(1,2)`) and trigrams added expressive power (e.g., capturing phrases like "not spam"), but it also significantly increased the feature count. I found a trade-off at `max_features` around 8000, which preserved most of the gains while keeping the model size manageable.
- **Insights & Understanding:**
 - **Baseline Importance:** Bag-of-Words (BoW) methods remain a very strong baseline. Despite the rise of complex transformer models, for small to medium-sized datasets, a well-tuned `CountVectorizer` or `TfidfVectorizer` combined with a simple classifier often matches or even outperforms deep models, and at a fraction of the computational cost.
 - **Interpretability:** BoW models preserve interpretability. By inspecting the highest-weighted tokens from the logistic regression model (via the `.coef_` attribute), I could immediately identify influential words (e.g., "free," "click," "unsubscribe" in the spam context).
- **Challenges Encountered:**
 - **Dimensionality & Sparsity:** My initial document-term matrix was very large (20,000 documents \times 100,000 terms), which strained memory and slowed down matrix multiplications. I pruned the vocabulary using `min_df=5` and consistently used the sparse CSR (Compressed Sparse Row) format to manage this.
 - **Overfitting with High-Order N-grams:** Adding trigrams without adequate document frequency filtering caused the classifier to overfit, leading to a 5% drop in validation accuracy. I introduced `max_df=0.8` to exclude extremely common n-grams, which helped.
- **Application & Relevance:**
 - **Scalable Text Pipelines:** I packaged the vectorizer and classifier into a scikit-learn Pipeline, which enables quick deployment and reproducible experiments.
 - **Industry Use Cases:** In past internships, I've built topic classification services for support tickets using similar BoW pipelines. This lab clarified parameter tuning strategies that I can now automate using techniques like grid search.
- **Code & Experimentation:**
 - **Grid Search:** I ran `GridSearchCV` over parameters like `max_features` ([2000, 5000, 8000]), `ngram_range` ([(1,1), (1,2)]), and `binary` ([True, False]). The best configuration I found was TF-IDF with (1,2) n-grams and `max_features=5000`.
 - **Alternative Models:** Beyond logistic regression, I also tried Support Vector Machines (SVM, specifically `LinearSVC`) and Naive Bayes. SVM matched logistic regression's accuracy (around 88%) but its training times were about twice as slow.

Lab 3: Using GloVe Embeddings

- **Summary of Key Learnings:**
 - **GloVe Embedding Concepts:** I learned that GloVe (Global Vectors for Word Representation) learns word vectors by factorizing a global word co-occurrence matrix. I gained a theoretical understanding from the original paper by Pennington et al. (2014) and saw how these co-occurrence probabilities translate into vector dot-product relationships.

- **Loading & Efficient Storage:** I parsed the large glove.840B.300d.txt file (822MB) into a Python dictionary mapping words to NumPy arrays. To speed up subsequent runs, I serialized this dictionary into a glove.npz file using `np.savez_compressed`, which reduced the loading time from about 2 minutes to just 5 seconds.
- **Cosine Similarity & Analogies:** Using simple cosine similarity, I verified classic word analogies like $\text{vector('king')} - \text{vector('man')} + \text{vector('woman')} \approx \text{vector('queen')}$. Quantitatively, the top-1 analogy accuracy on a small evaluation set I used was around 72%.
- **Insights & Understanding:**
 - **Semantic Richness:** Embeddings truly capture analogical structure. Visualizing the top two principal components (via PCA) of word vectors revealed distinct clusters for days of the week, countries, and vehicle terms, showcasing their semantic richness.
 - **Transfer Learning:** Pre-trained embeddings eliminate the need to learn word representations from scratch, especially when data is limited. This significantly reduces training time and data requirements for downstream NLP tasks.
- **Challenges Encountered:**
 - **OOV (Out-of-Vocabulary) Handling:** About 15% of the tokens in my downstream sentiment analysis dataset were not present in the GloVe vocabulary. I implemented three strategies to handle these OOV words: random initialization, using a zero vector, and averaging the vectors of known neighboring words. Averaging neighbors gave slightly better validation accuracy (a +1% improvement).
 - **Normalization Effects:** In my early experiments, not normalizing the GloVe vectors (i.e., using their raw magnitudes) led to unstable training. Normalizing each vector to unit length before feeding it into the model stabilized the gradients.
- **Application & Relevance:**
 - **Integration into Neural Models:** I initialized a Keras Embedding layer with GloVe weights. I first kept these weights frozen and then experimented with fine-tuning them in a later stage. Fine-tuning improved the test accuracy by about 2%.
 - **Beyond NLP:** I realized that the concept of embedding high-dimensional categorical variables extends beyond NLP to areas like recommender systems and graph embeddings.
- **Code & Experimentation:**
 - **PCA Visualization:** I wrote a quick Matplotlib script to plot 500 selected word vectors in 2D, which clearly revealed these semantic clusters.
 - **Analogy Benchmark:** I automated the evaluation on the Google analogy test set, achieving around 68% on semantic analogies and 55% on syntactic analogies.

Lab 4: Introducing RNNs

- **Summary of Key Learnings:**
 - **Sequence Preparation:** I tokenized text using Keras's Tokenizer (setting `num_words=10000`), converted the text to sequences of integers, and then applied `pad_sequences` (with `maxlen=100`). I learned that using shorter sequences (e.g., ≤ 100 tokens) can speed up training while still capturing key context.
 - **Embedding + LSTM:** I built an `nn.Embedding` layer (with `vocab_size=10000` and `embedding_dim=100`) and initialized it with frozen GloVe weights. I then added an

nn.LSTM layer (with hidden_size=128, num_layers=1, and batch_first=True), followed by two nn.Linear layers for the classification task.

- **Training Loop & Metrics:** I trained the model for 10 epochs with a learning rate of $1e-3$, a batch size of 32, using the AdamW optimizer. I tracked training and validation loss and accuracy at each epoch, achieving about 82% validation accuracy on a binary sentiment task.
- **Insights & Understanding:**
 - **Vanishing/Exploding Gradients:** I noted that using longer sequences (e.g., length 200) led to gradient explosion. Applying gradient clipping (`torch.nn.utils.clip_grad_norm_`) with a threshold of 5 stabilized the training process.
 - **Freezing vs. Fine-Tuning Embeddings:** Freezing the embedding layer during the initial epochs helped prevent "catastrophic forgetting" of the pre-trained GloVe knowledge. Later unfreezing the embedding layer for fine-tuning gave a final accuracy boost of about 2%.
- **Challenges Encountered:**
 - **Batch Size vs. Sequence Length Trade-off:** On my GPU (with 8GB of memory), using a batch size of 32 with a sequence length of 200 caused out-of-memory errors. I settled on a sequence length of 100 and a batch size of 32 to manage this.
 - **Training Instability:** The initial training loss sometimes spiked unexpectedly. Adding LayerNorm after the LSTM output and before the classifier layer helped stabilize the learning process and improved convergence speed by about 20%.
- **Application & Relevance:**
 - **Extensibility to Attention Models:** Recognizing the limitations of RNNs in handling long-range dependencies, I sketched out plans to add self-attention layers or even transition to a Transformer encoder architecture in future work.
 - **Time-Series Applications:** The skills I gained in sequence modeling are directly transferable to forecasting tasks, such as analyzing sensor data or stock prices.
- **Code & Experimentation:**
 - **Bidirectional LSTM:** I tested setting `bidirectional=True` on the LSTM layer. This improved validation accuracy from 82% to 85% but also doubled the training time.
 - **Early Stopping & Checkpointing:** I implemented early stopping with a patience of 3 (monitoring validation loss) and saved the best model checkpoint to disk. This prevented overfitting, which started to occur after epoch 7 in one of my runs.

Lab-06: AWS Machine Learning University Module 2 Lab 05 - Fine Tuning BERT

1. Learning Insights:

- **1.1 Transfer Learning with Transformers:** Instead of training a language model from the ground up, I loaded the `bert-base-uncased` model using Hugging Face's Transformers library. I then repurposed its deep bidirectional encodings for a sentiment classification task. The fine-tuning process itself only involved adding and training a lightweight

classification head on top of BERT's pooled output. It was amazing to see how much contextual knowledge was already embedded in the pretrained weights, derived from over 3 billion words from Wikipedia and BookCorpus.

- **1.2 Tokenization and Input Formatting:** I gained a much deeper appreciation for subword tokenization, specifically WordPiece, which BERT uses. Converting raw text into `input_ids`, `attention_masks`, and `token_type_ids` taught me how BERT handles out-of-vocabulary words by breaking them down into known subwords (e.g., "unbelievable" becomes "un" and "##believable"). I also learned the importance of padding sequences to a fixed `max_length` (we used 128) and batching examples efficiently to ensure the model receives uniformly shaped tensors.
- **1.3 Optimization & Learning-Rate Scheduling:** The lab guided me to use the AdamW optimizer with a very low learning rate (2e-5) and a linear warmup schedule for the first 10% of training steps. Observing the training curves clearly showed how this warmup phase prevents large initial weight updates that could otherwise destroy the valuable pretrained features. It highlighted that even state-of-the-art models require careful tuning of optimizer hyperparameters for smooth convergence.
- **1.4 Evaluation Metrics:** Beyond just looking at accuracy, I computed precision, recall, and the F1-score on the test set. The fine-tuned BERT model achieved around 92% accuracy and an F1-score of approximately 0.91 on a balanced reviews dataset. Tracking these different metrics reinforced the idea that accuracy alone can sometimes mask issues like class imbalance, and that the F1-score is often a better indicator of real-world performance for binary classification tasks.

2. Challenges and Struggles:

- **2.1 Resource Constraints & Batch Size:** Fine-tuning BERT on my consumer-grade GPU (8GB) quickly led to Out-of-Memory (OOM) errors when I tried using batch sizes greater than 16. I experimented with gradient accumulation, where I accumulated gradients over 4 steps with a batch size of 8 to effectively emulate a batch size of 32. This allowed for stable training without running out of memory.
- **2.2 Tokenizer Quirks:** When I first created the dataset, I passed raw reviews directly to BERT's `tokenizer.encode_plus` method. I initially forgot to set `truncation=True`, so some reviews exceeded the model's maximum length, causing index errors during processing. Adding `truncation=True` and `padding='max_length'` fixed this issue.
- **2.3 Overfitting & Early Stopping:** After about 4 epochs, I noticed the training loss continued to decrease, but the validation loss started to rise – a classic sign of overfitting. Implementing early stopping (with a patience of 2, monitoring validation loss) and saving the best checkpoint prevented the model's performance from degrading further.
- **2.4 Understanding Model Outputs:** Initially, I mistakenly treated BERT's raw output logits as probabilities. It was only after applying a softmax function that the outputs summed to 1.0 and became properly interpretable as probabilities. Plotting the softmax outputs for a few examples helped me debug why the model was making some confidently wrong predictions.

3. Personal Growth:

- **3.1 Confidence with Hugging Face APIs:** Before this lab, I had mostly read about Transformers but hadn't used them extensively hands-on. Now, I feel much more comfortable instantiating BertTokenizer, BertForSequenceClassification, and defining a PyTorch Dataset that correctly feeds input_ids, attention_mask, and labels into the model.
- **3.2 Pipelines for Reproducibility:** I made an effort to refactor the training and evaluation steps into reusable functions (like train_epoch() and eval_epoch()). I also wrote a simple command-line argument parser so I could rerun experiments with different hyperparameters without needing to rewrite code. Adopting this "production mindset" was a new and incredibly valuable experience for me.
- **3.3 Broader Perspective on NLP:** Fine-tuning BERT really drove home the point that modern NLP workflows heavily revolve around using large pretrained models and then adapting them to specific domains or tasks. I now see how these techniques extend far beyond sentiment analysis to areas like question answering, named-entity recognition, and more, opening up many exciting avenues for future projects.

4. Critical Reflection:

- **4.1 What I Would Do Differently:**
 - **Mixed-Precision Training:** I would integrate NVIDIA's AMP (Automatic Mixed Precision) to potentially reduce memory usage further and speed up training.
 - **Dynamic Padding:** Instead of padding all sequences in a batch to the same maximum length, I would explore sorting examples by length within each batch to minimize the number of wasted padding tokens.
 - **Cross-Validation:** For a more robust estimate of generalization performance, especially on smaller datasets, I would implement k-fold cross-validation.
- **4.2 New Questions & Exploration:**
 - **Model Size vs. Speed:** How much accuracy would I realistically lose if I opted for a smaller, faster model like DistilBERT or ALBERT?
 - **Domain Adaptation:** Could further pretraining on highly specific in-domain text (e.g., logs from a Mars rover) improve performance for a niche task?
 - **Explainability:** What methods, like attention-based heatmaps, can I apply to better interpret which words or tokens most influence BERT's decisions?
- **4.3 Fit into the Broader ML Landscape:** This lab cemented my understanding that transfer learning with large-scale pretrained models is truly at the current frontier of NLP. It complements earlier work with sequence-based models like RNNs and clearly points the way towards fully self-supervised, task-agnostic representations—the path that has led to models like GPT and other Large Language Models (LLMs).

