**KENT STATE**
U N I V E R S I T Y

# Assignment 4: Text and Sequence Data

## 1. Executive Summary

The fast development of natural language processing (NLP) has raised the interest of deep learning architectures which are able to interpret sequential data and especially text. Recurrent Neural Networks (RNNs) and Transformer-based constructions, among others, have been used as a fundamental component in such activities as sentiment analysis, translation, and classification. They are successful not only in terms of architecture but also in the ways of textual input preparation, representation and integration in the process of learning. These considerations are discussed in this assignment where the IMDB sentiment analysis dataset is analyzed in the deliberately constrained settings with limited training size, vocabulary, and shortened review lengths.

This assignment is aimed at testing the effectiveness of the application of RNNs or Transformers to text data and to learn the behavior of embedding-based strategies in different training data conditions. In particular, the experiment contrasts learned embeddings, those that are directly trained on the sentiment categorization training, with pretrained embeddings, which also give general information on the semantics but do not change during the training process. The assignment can detect the performance trend, strength, and weakness of the various embedding options by observing their model performance throughout a variety of training sample sizes.

The main theme of the assignment is how the models that have been trained with little supervision can be useful in meaningful predictions and what embedding strategy is the most appropriate to enhance the results of the prediction. This exploration provides an understanding of model architecture as well as a concept of data efficiency, which is a vital attribute of practical machine learning issues, in which large amounts of labeled data are frequently unattainable. The assignment offers a practical insight into the trade-offs between the pretraining and task-specific learning of text-based deep learning through a structured experimentation.

## 2.  Application of RNNs or Transformers to a Text and Sequence Data.

To use RNNs or Transformers on text data, one will have to convert raw text into a standard, numerical format, which neural networks will work with. This starts by tokenization where the reviews are divided into separate words and the words are transformed into distinct integer values. The text is restricted to the top 10,000 most common words to obtain stable calculations, and the length of each individual review is shortened or extended to a fixed length of 150 tokens. The steps ensure that the model gets equal length numerical sequences no matter what the input may vary. Having such representation, the model can understand text effectively and analyze it.

The second step, which follows tokenization, is embedding that is an indexing of word indices into dense vector representations. This may include trained embeddings which are initialized at random and updated through training or trained embeddings with contextual significance which may be an external corpus like GloVe that are rich and provide deep meaning. The layer of embedding retrieves semantic relations and transforms discrete tokens into meaningful vectors that indicate word similarities. After embedding, the sequence may be input into either an RNN-based architecture such as LSTM which learns the time-dependence of the sequence by processing tokens sequentially, or a Transformer encoder which learns the relationships between contextual information globally by using self-attention on the entire sequence.

Lastly, an embedded and processed sequence is sent to a neural network classifier. Regularization also acts to avoid overfitting like dropout, when there is a small amount of training data. A sigmoid-shaped thick layer would give binary sentiment prediction. With this pipeline, being tokenized, embedded, sequentially modelled, and classified, RNNs and Transformers are

successfully trained on the IMDB dataset. This assignment follows these principles to test the effects of embedding choice on performance in a controlled setting.

### 3. Dealing with limited data on How to Improve Performance.

One of the most important problems of deep learning is how to improve the performance of a model with limited data. In case of the very limited sample size of training as in this assignment, which is 100, the model does not have enough data to learn any powerful patterns directly out of the data. One of the best answers is the introduction of pretrained embeddings, which encode semantic meaning, having accessed millions of external text in the past. The model uses these useful representations and therefore avoids having to re-learn the basic structure of the language and attains more stable performance with a limited amount of data.

Techniques of regularization are also important. The use of dropout layers compels the model to follow more than one path to prediction rather than memorizing some particular pattern, which is specific to the small dataset, thus enhancing generalization. Early termination helps the model to avoid training past the stage of improvement which reduces the risk of overfitting. Restricting the maximum words of the vocabulary to 10,000 and word truncated reviews to 150 words also contribute to it because they simplify the input space and make sure the model only learns the most significant information. These preprocessing mechanisms are architectural protection mechanisms that assist the model to be more effective even when the data are severely limited.

KENT STATE
U N I V E R S I T Y

Naturally, performance also increases with the size of training. Once the dataset is increased to 200, 500, 1 000 or more training samples, the learned embedding model starts to learn subtle sentiment patterns that cannot be learnt by the pretrained embeddings. Knowledge of the effect of sample size on performance can be used to determine the point at which learned embeddings can be more effective (crossover point). Pretrained embeddings, regularization methods, and well-planned preprocessing make a combination of these components to increase performance in the conditions of low resources.

### 4. Identifying the most appropriate methods to be used in prediction improvement.

The assignment makes a direct parallelogram of learnt embeddings and pretrained embeddings to identify the method that proves to be more appropriate regarding prediction accuracy. In the situation where there are few labeled examples, pretrained embeddings performed the most significant improvement since they imply intrinsic linguistic mechanism. This is proven by the graph uploaded: at the minimal training sizes (approximately 100 samples), all three models begin at approximately 0.49-0.50 accuracy, with pretrained frozen and fine-tuned models being slightly better than learned embeddings.

The performance trend was however reversed as the training samples increased. At 300-500 samples the learned embeddings (blue line) start to increase much more quickly than either of the pretrained models, attaining about 0.79 accuracy with frozen pretrained embeddings at about 0.53 and fine-tuned embeddings at about 0.57. Such change is an indication of the ability of the learned embeddings to adjust to IMDB-specific sentiment patterns, such as domain-specific language and contextual indicators. In the meantime, it is not possible to add dataset-specific nuances to pretrained embeddings, and in
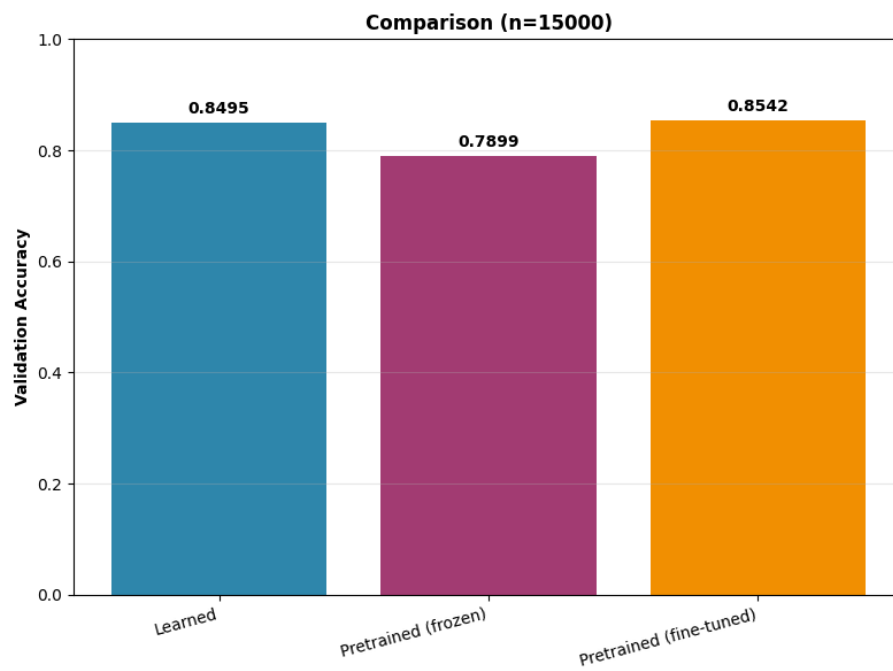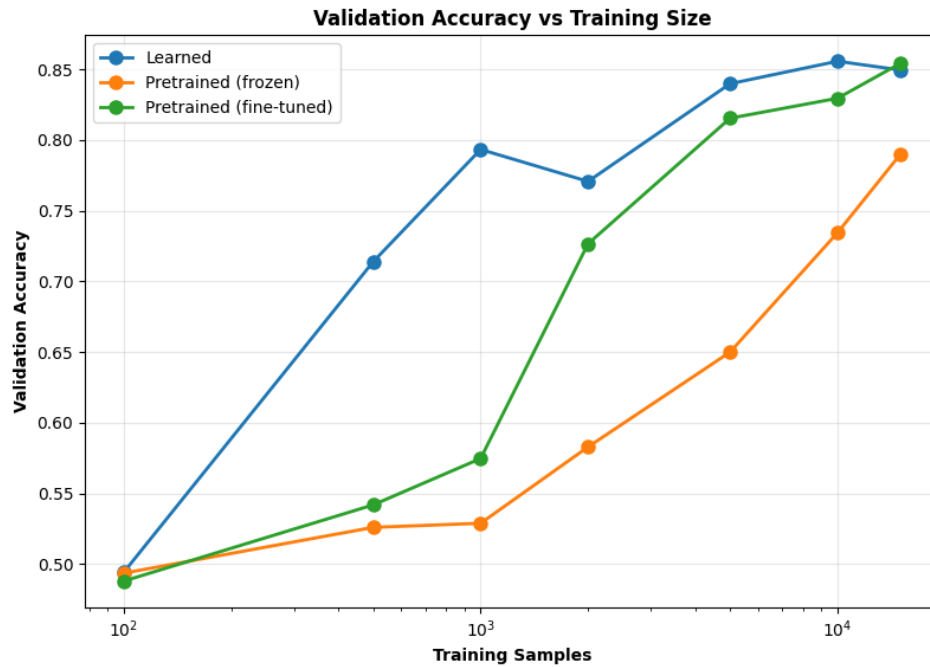
frozen ones.

Finally, the appropriateness is based on the size of dataset. Embeddings that are pretrained are better to enhance performance in cases where the data available is very scarce but learned embeddings are more appropriate in cases where information is sufficient to allow the model to acquire task-specific patterns. Comparison bar chart at large sample sizes, especially in the range of 15,000, the comparison bar chart reveals the fine-tuned pretrained embeddings to be best at 0.8542 and the learned embeddings at 0.8495, with the frozen pretrained embeddings coming in last at 0.7899. This shows the great concomitant power of pretraining and fine-tuning when sufficient data is present.

### 5. Table of Results

## Validation Accuracy by Training Sample Size

| Training Samples | Learned Embeddings | Pretrained (Fine-tuned) | Pretrained (Frozen) |
|---|---|---|---|
| 100 | 0.4942 | 0.4880 | 0.4936 |
| 500 | 0.7139 | 0.5419 | 0.5260 |
| 1,000 | 0.7933 | 0.5747 | 0.5289 |
| 2,000 | 0.7707 | 0.7266 | 0.5829 |
| 5,000 | 0.8398 | 0.8154 | 0.6500 |
| 10,000 | 0.8557 | 0.8295 | 0.7344 |
| 15,000 | 0.8495 | 0.8542 | 0.7899 |

Validation Accuracy vs Training Size



Comparison (n=15000)

KENT STATE
U N I V E R S I T Y

## 6. Results and Analysis Discussion.

The findings show clearly the change in model performance with the change in the size of the training and the embedding strategy employed. Under the smallest data regimes of approximately 100 training samples, all three models start at approximately 0.49-0.50 accuracy with pretrained embeddings doing slightly better than learned embeddings. This holds to the fact that pretrained embeddings are more adapted to be used in highly under-supervised environments where they offer a rich semantic prior that the model can exploit off-the-shelf.

The turning point is reached when the size of training reaches the mid-range (300-1000 samples), and learned embeddings start to outperform not only frozen pretrained embeddings but also fine-tuned pretrained embeddings. The learned model quickly rises above 0.70 and it becomes nearly 0.80 accurate after approximately 1,000 samples. Contrarily though, frozen pretrained embeddings are close to 0.52-0.53, and fine-tuned embeddings are behind until approximately 2,000 samples. This shows that learned embeddings learn faster towards the task-specific IMDB sentiment structure.

Models are improved with increasing size of training (5,000-15,000 samples) but the ranking changes once again. The best accuracy (0.8542) is finally achieved at 15,000 samples. The learned embeddings are highly competitive with a value of 0.8495, and the frozen pretrained embeddings are lowest. This last trend shows that pretrained semantic knowledge used together with the adaptability of fine-tuning is the most effective on a large scale.

**KENT STATE**
U N I V E R S I T Y

## 7. Conclusion

The assignment was a detailed analysis of the application of RNNs and Transformer-based models to text data, the possibility of getting better performance in limited training settings, and embedding strategies that present the highest performance in terms of prediction. Constrained to work with the IMDB data, the experiment identified the strengths and weaknesses of learned, pretrained frozen, and pretrained fine-tuned embeddings. The pretrained embeddings were useful in low-data settings, where the semantic priors of the embeddings ensured a stabilized performance level. Nevertheless, learned embeddings were better with the training set size of a few hundred samples, with learned embeddings outperforming pretrained frozen embeddings as well as fine-tuned ones when the training set size was in the mid-range regime.

The findings highlight a key concept in machine learning, which is that a model is only as good as its architecture and the accessibility of data and appropriateness of embedding strategies. The power of general semantic knowledge and task-specific learning can be seen at very large data sizes where the fine-tuned pretrained model provided the most effective overall accuracy. In the end, the assignment effectively demonstrates the relationship between data size, embedding strategy and model performance, which gives credence to the basic principles of text-based deep learning.

## 8. Appendix

# Assignment 4: Text and Sequence Data - IMDB Sentiment Analysis

**Author:** Divya Ratakonda

## 1. Setup and Imports

```python
import os
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
from pathlib import Path
import glob


np.random.seed(42)


import warnings
warnings.filterwarnings('ignore')


os.environ["KERAS_BACKEND"] = "tensorflow"

import keras
from keras import layers
from keras.datasets import imdb
from keras.utils import pad_sequences
```

## 2. Configuration

```python
MAX_FEATURES = 10000    # Top 10,000 words
MAX_LEN = 150           # Cutoff reviews after 150 words
VAL_SAMPLES = 10000     # Validate on 10,000 samples
EMBEDDING_DIM = 100     # GloVe 6B: 50, 100, 200, 300 (NO 128!)
HIDDEN_DIM = 64
EPOCHS = 20
BATCH_SIZE = 32

TRAINING_SAMPLE_SIZES = [100, 500, 1000, 2000, 5000, 10000, 15000, 20000]
```

## 3. IMDB Dataset

```python
# Load IMDB data
(x_train_full, y_train_full), (x_test, y_test) = imdb.load_data(num_words=MAX_FEATURES)

print(f"Full training: {len(x_train_full)} samples")
print(f"Test: {len(x_test)} samples")

# Pad sequences
x_train_full = pad_sequences(x_train_full, maxlen=MAX_LEN)
x_test = pad_sequences(x_test, maxlen=MAX_LEN)

# Split
```

```
x_val = x_train_full[:VAL_SAMPLES]
y_val = y_train_full[:VAL_SAMPLES]
x_train_available = x_train_full[VAL_SAMPLES:]
y_train_available = y_train_full[VAL_SAMPLES:]

print(f"\nAfter split:")
print(f"  Training available: {len(x_train_available)}")
print(f"  Validation: {len(x_val)}")
```

```
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/imdb.npz
17464789/17464789 ─────────────── 0s 0us/step
Full training: 25000 samples
Test: 25000 samples

After split:
  Training available: 15000
  Validation: 10000
```

## 4. Load GloVe Embeddings

```
def find_glove_file(embedding_dim):

    filename = f"glove.6B.{embedding_dim}d.txt"


    search_paths = [
        f"./{filename}",
        f"./glove/{filename}",
        Path.home() / f".keras/datasets/{filename}",
        Path.home() / f".keras/datasets/glove/{filename}",
    ]

    print(f"Searching for {filename}...")


    for path in search_paths:
        path = Path(path)
        if path.exists():
            print(f" Found at: {path}")
            return path


    print("Searching recursively...")
    glob_patterns = [
        f"**/{filename}",
    ]

    for pattern in glob_patterns:
        matches = glob.glob(pattern, recursive=True)
        if matches:
            print(f" Found at: {matches[0]}")
            return Path(matches[0])

    return None
```

```
def load_glove_embeddings(embedding_dim, max_words):


    print(f"Attempting to download/locate GloVe {embedding_dim}d...")
    try:
        glove_url = "http://nlp.stanford.edu/data/glove.6B.zip"
        zip_path = keras.utils.get_file(
            "glove.6B.zip",
            origin=glove_url,
```

```
            extract=True,
            cache_dir='.',
            cache_subdir='glove'
        )


        glove_dir = Path(zip_path).parent
        glove_path = glove_dir / f"glove.6B.{embedding_dim}d.txt"

        if glove_path.exists():
            print(f" Found at: {glove_path}")
        else:
            print("Not in expected location, searching...")
            glove_path = find_glove_file(embedding_dim)

            if glove_path is None:
                raise FileNotFoundError(f"Cannot find glove.6B.{embedding_dim}d.txt")

    except Exception as e:
        print(f"Info: {e}")
        print("Searching for existing file...")
        glove_path = find_glove_file(embedding_dim)

        if glove_path is None:
            raise FileNotFoundError(
                f"Could not find glove.6B.{embedding_dim}d.txt\n"
                f"Download from: http://nlp.stanford.edu/data/glove.6B.zip\n"
                f"Extract and place in current directory"
            )

    # Load embeddings
    print(f"Loading embeddings from: {glove_path}")
    embeddings_index = {}

    with open(glove_path, 'r', encoding='utf-8') as f:
        for line in f:
            values = line.split()
            word = values[0]
            coefs = np.asarray(values[1:], dtype='float32')
            embeddings_index[word] = coefs

    print(f"Loaded {len(embeddings_index)} word vectors")

    # Create embedding matrix
    embedding_matrix = np.zeros((max_words, embedding_dim))
    word_index_imdb = imdb.get_word_index()
    reverse_word_index = {v: k for k, v in word_index_imdb.items()}

    found_words = 0
    for i in range(1, max_words):
        word = reverse_word_index.get(i)
        if word is not None:
            embedding_vector = embeddings_index.get(word)
            if embedding_vector is not None:
                embedding_matrix[i] = embedding_vector
                found_words += 1

    print(f" Mapped {found_words}/{max_words} words ({100*found_words/max_words:.1f}%)")

    return embedding_matrix

print(" GloVe  function defined")

 GloVe  function defined
```

```
try:
    embedding_matrix = load_glove_embeddings(EMBEDDING_DIM, MAX_FEATURES)
    glove_available = True
```

```
        print("\n GloVe embeddings loaded")
    except Exception as e:
        print(f"\n Could not load GloVe: {e}")
        print("Continuing with learned embeddings only")
        glove_available = False
```

```
Attempting to download/locate GloVe 100d...
Downloading data from http://nlp.stanford.edu/data/glove.6B.zip
862182613/862182613 ━━━━━━━━━━━━━━━ 161s 0us/step
Not in expected location, searching...
Searching for glove.6B.100d.txt...
Searching recursively...
 Found at: glove/glove_extracted/glove.6B.100d.txt
Loading embeddings from: glove/glove_extracted/glove.6B.100d.txt
Loaded 400000 word vectors
Downloading data from https://storage.googleapis.com/tensorflow/tf-keras-datasets/imdb_word_index.json
1641221/1641221 ━━━━━━━━━━━━━━━ 0s 0us/step
 Mapped 9796/10000 words (98.0%)

 GloVe embeddings loaded
```

## 5. Define Model Architectures

```python
def build_lstm_learned(max_features, embedding_dim, max_len, hidden_dim):

    inputs = keras.Input(shape=(max_len,), dtype="int32")
    x = layers.Embedding(max_features, embedding_dim, mask_zero=True)(inputs)
    x = layers.Bidirectional(layers.LSTM(hidden_dim))(x)
    x = layers.Dropout(0.5)(x)
    outputs = layers.Dense(1, activation="sigmoid")(x)

    model = keras.Model(inputs, outputs, name="lstm_learned")
    model.compile(optimizer="adam", loss="binary_crossentropy", metrics=["accuracy"])
    return model

def build_lstm_pretrained(embedding_matrix, max_len, hidden_dim, trainable=False):
    """LSTM with pretrained embeddings."""
    max_features, embedding_dim = embedding_matrix.shape

    inputs = keras.Input(shape=(max_len,), dtype="int32")
    x = layers.Embedding(
        max_features, embedding_dim,
        weights=[embedding_matrix],
        trainable=trainable,
        mask_zero=True
    )(inputs)
    x = layers.Bidirectional(layers.LSTM(hidden_dim))(x)
    x = layers.Dropout(0.5)(x)
    outputs = layers.Dense(1, activation="sigmoid")(x)

    model = keras.Model(inputs, outputs, name="lstm_pretrained")
    model.compile(optimizer="adam", loss="binary_crossentropy", metrics=["accuracy"])
    return model

print(" Model architectures defined")
```

```
 Model architectures defined
```

## 6. Training Loop

```python
results = []
early_stopping = keras.callbacks.EarlyStopping(
    monitor="val_loss", patience=3, restore_best_weights=True
)
```

```
    for n_samples in TRAINING_SAMPLE_SIZES:

        print(f"Training with {n_samples} samples")


        if n_samples > len(x_train_available):
            continue

        x_train = x_train_available[:n_samples]
        y_train = y_train_available[:n_samples]

        # Model 1: Learned
        print("\n[1/3] Learned embeddings...")
        model = build_lstm_learned(MAX_FEATURES, EMBEDDING_DIM, MAX_LEN, HIDDEN_DIM)
        model.fit(x_train, y_train, batch_size=BATCH_SIZE, epochs=EPOCHS,
                  validation_data=(x_val, y_val), callbacks=[early_stopping], verbose=0)
        _, val_acc = model.evaluate(x_val, y_val, verbose=0)
        print(f" Val Acc: {val_acc:.4f}")
        results.append({'n_samples': n_samples, 'model': 'Learned', 'val_accuracy': val_acc})

        if glove_available:
            # Model 2: Pretrained (frozen)
            print("\n[2/3] Pretrained (frozen)...")
            model = build_lstm_pretrained(embedding_matrix, MAX_LEN, HIDDEN_DIM, trainable=False)
            model.fit(x_train, y_train, batch_size=BATCH_SIZE, epochs=EPOCHS,
                      validation_data=(x_val, y_val), callbacks=[early_stopping], verbose=0)
            _, val_acc = model.evaluate(x_val, y_val, verbose=0)
            print(f" Val Acc: {val_acc:.4f}")
            results.append({'n_samples': n_samples, 'model': 'Pretrained (frozen)', 'val_accuracy': val_a

            # Model 3: Pretrained (fine-tuned)
            print("\n[3/3] Pretrained (fine-tuned)...")
            model = build_lstm_pretrained(embedding_matrix, MAX_LEN, HIDDEN_DIM, trainable=True)
            model.fit(x_train, y_train, batch_size=BATCH_SIZE, epochs=EPOCHS,
                      validation_data=(x_val, y_val), callbacks=[early_stopping], verbose=0)
            _, val_acc = model.evaluate(x_val, y_val, verbose=0)
            print(f" Val Acc: {val_acc:.4f}")
            results.append({'n_samples': n_samples, 'model': 'Pretrained (fine-tuned)', 'val_accuracy': \

    print("\n Training complete!")
```

```
Training with 100 samples

[1/3] Learned embeddings...
 Val Acc: 0.4942

[2/3] Pretrained (frozen)...
 Val Acc: 0.4936

[3/3] Pretrained (fine-tuned)...
 Val Acc: 0.4880
Training with 500 samples

[1/3] Learned embeddings...
 Val Acc: 0.7139

[2/3] Pretrained (frozen)...
 Val Acc: 0.5260

[3/3] Pretrained (fine-tuned)...
 Val Acc: 0.5419
Training with 1000 samples

[1/3] Learned embeddings...
 Val Acc: 0.7933

[2/3] Pretrained (frozen)...
 Val Acc: 0.5289

[3/3] Pretrained (fine-tuned)...
 Val Acc: 0.5747
Training with 2000 samples
```

```
[1/3] Learned embeddings...
 Val Acc: 0.7707

[2/3] Pretrained (frozen)...
 Val Acc: 0.5829

[3/3] Pretrained (fine-tuned)...
 Val Acc: 0.7266
Training with 5000 samples

[1/3] Learned embeddings...
 Val Acc: 0.8398

[2/3] Pretrained (frozen)...
 Val Acc: 0.6500

[3/3] Pretrained (fine-tuned)...
 Val Acc: 0.8154
Training with 10000 samples

[1/3] Learned embeddings...
 Val Acc: 0.8557

[2/3] Pretrained (frozen)...
 Val Acc: 0.7344
```

## 7. Results Analysis

```python
results_df = pd.DataFrame(results)

print("Detailed Results:")
print(results_df.to_string(index=False))

if len(results_df) > 0:
    pivot = results_df.pivot(index='n_samples', columns='model', values='val_accuracy')
    print("\n\nValidation Accuracy by Training Sample Size:")
    print(pivot.to_string())
```

```
Detailed Results:
 n_samples                 model  val_accuracy
       100              Learned        0.4942
       100      Pretrained (frozen)   0.4936
       100 Pretrained (fine-tuned)    0.4880
       500              Learned        0.7139
       500      Pretrained (frozen)   0.5260
       500 Pretrained (fine-tuned)    0.5419
      1000              Learned        0.7933
      1000      Pretrained (frozen)   0.5289
      1000 Pretrained (fine-tuned)    0.5747
      2000              Learned        0.7707
      2000      Pretrained (frozen)   0.5829
      2000 Pretrained (fine-tuned)    0.7266
      5000              Learned        0.8398
      5000      Pretrained (frozen)   0.6500
      5000 Pretrained (fine-tuned)    0.8154
     10000              Learned        0.8557
     10000      Pretrained (frozen)   0.7344
     10000 Pretrained (fine-tuned)    0.8295
     15000              Learned        0.8495
     15000      Pretrained (frozen)   0.7899
     15000 Pretrained (fine-tuned)    0.8542


Validation Accuracy by Training Sample Size:
model      Learned  Pretrained (fine-tuned)  Pretrained (frozen)
n_samples
100         0.4942                   0.4880               0.4936
500         0.7139                   0.5419               0.5260
1000        0.7933                   0.5747               0.5289
2000        0.7707                   0.7266               0.5829
5000        0.8398                   0.8154               0.6500
```

|       |        |        |        |
|-------|--------|--------|--------|
| 10000 | 0.8557 | 0.8295 | 0.7344 |
| 15000 | 0.8495 | 0.8542 | 0.7899 |

## 8. Visualization

```python
fig, axes = plt.subplots(1, 2, figsize=(16, 6))

# Plot 1: Line plot
ax1 = axes[0]
for model_name in results_df['model'].unique():
    data = results_df[results_df['model'] == model_name]
    ax1.plot(data['n_samples'], data['val_accuracy'],
             marker='o', linewidth=2, markersize=8, label=model_name)

ax1.set_xlabel('Training Samples', fontweight='bold')
ax1.set_ylabel('Validation Accuracy', fontweight='bold')
ax1.set_title('Validation Accuracy vs Training Size', fontweight='bold')
ax1.legend()
ax1.grid(True, alpha=0.3)
ax1.set_xscale('log')

# Plot 2: Bar chart
ax2 = axes[1]
if len(results_df) > 0:
    latest = results_df[results_df['n_samples'] == results_df['n_samples'].max()]
    bars = ax2.bar(range(len(latest)), latest['val_accuracy'],
                   color=['#2E86AB', '#A23B72', '#F18F01'])
    ax2.set_ylabel('Validation Accuracy', fontweight='bold')
    ax2.set_title(f'Comparison (n={int(latest["n_samples"].max())})', fontweight='bold')
    ax2.set_xticks(range(len(latest)))
    ax2.set_xticklabels(latest['model'], rotation=15, ha='right')
    ax2.set_ylim([0, 1])
    ax2.grid(True, alpha=0.3, axis='y')

    for bar, acc in zip(bars, latest['val_accuracy']):
        ax2.text(bar.get_x() + bar.get_width()/2, bar.get_height() + 0.01,
                 f'{acc:.4f}', ha='center', va='bottom', fontweight='bold')

plt.tight_layout()
plt.savefig('assignment4_results.png', dpi=300, bbox_inches='tight')
print("✓ Saved: assignment4_results.png")
plt.show()
```

aved: assignment4_results.png



**Validation Accuracy vs Training Size**

```python
# Save CSV
results_df.to_csv('assignment4_results.csv', index=False)
print(" Saved: assignment4_results.csv")

# Print conclusions

print("CONCLUSIONS")


if len(results_df) > 0:
    best = results_df.loc[results_df['val_accuracy'].idxmax()]
    print(f"\n Best model: {best['model']}")
    print(f"  Accuracy: {best['val_accuracy']:.4f}")
    print(f"  Training samples: {int(best['n_samples'])}")

    small = results_df[results_df['n_samples'] <= 500]
    if len(small) > 0:
        best_small = small.loc[small['val_accuracy'].idxmax()]
        print(f"\nBest with ≤500 samples: {best_small['model']}")
        print(f"  Accuracy: {best_small['val_accuracy']:.4f}")
```

```
 Saved: assignment4_results.csv
CONCLUSIONS

 Best model: Learned
  Accuracy: 0.8557
  Training samples: 10000

Best with ≤500 samples: Learned
  Accuracy: 0.7139
```