

IMDB Sentiment Analysis: Using Neural Networks

1. Executive Summary

One of the most popular benchmarks in the natural language processing is the IMDB sentiment dataset. It houses 50,000 movie reviews which are equally split into positive and negative view. To ensure that a computer can use this information, each review is converted to a series of numbers, and each number is that of a word.

This assignment was not just aimed at creating a model that is capable of properly classifying reviews, but also to learn the impact of various design decisions on performance. It was possible to vary the architecture of the network by varying elements like the number of layers, units, and the training methods, and observe how the architecture influences the generalization capabilities of the network.

The assignment is directly related to two learning outcomes:

- 1) Acquiring hands-on experience to use Keras and TensorFlow to normalize and train a deep learning model.
- 2) To describe the effect of a particular modeling decision on accuracy and generalization. To put it in another way, it is equally important to know the reason why models act the way they do as much as it is getting the high accuracy.

2. How the Experiment was Setup

To prepare Reviews:

Only 10, 000 most frequent words of the dataset were retained to make the text manageable in the model. Individual reviews were then shrunk or fattened to precisely 200

words to make all of the individual reviews look consistent. This formed a uniform training format.

Words Being Thought to the Model:

A 128-number vector was obtained (so-called embedding) instead of working with raw text. The fixings were not predetermined but were learned in the course of the training so that a model could develop its own perception of the meaning of words gradually.

To Train the Model:

The common aspect of all versions of the model was the same training rules: it optimized Adam, a batch size of 512, and 20 percent validation split. The models were trained in 5 epochs, and some of them were trained in 8 epochs to examine the trends in greater detail. The last layer would always generate one probability, through a sigmoid function, which was an indicator of both positive and negative reviews.

What we Changed:

We varied various aspects of the model in a systematic manner in order to find out the effects of design choices on performance:

List of hidden Layers: 1,2,3

Units used: 32, 64, 128

Activation Functions used: tanh and Relu.

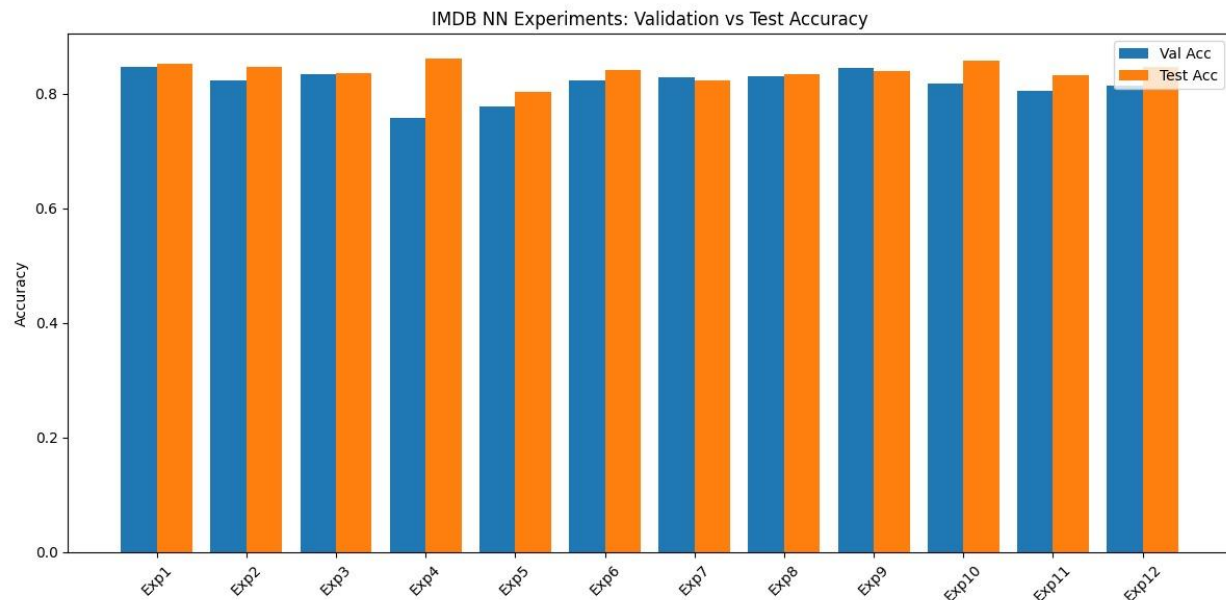
Loss Functions Binary Cross-Entropy (BCE) or Mean Squared Error (MSE)

Regularization: 0.3/ 0.5 dropout. Keeping some of the elements constant and varying the other, we might easily see the effects of each decision on the accuracy and generalization.

3. Results

Big picture: What were the most successful Models?

In the process of comparing all experiments to each other, it became obvious that models with two layers and 64 units always performed the strongest. Single-layer simplistic models fell slightly behind, and multidimensional networks were prone to overfitting, learning the training data too well that they would not be as effective at generalising.



Layers: Finding the Gold Middle Ground.

The extra layers added to the model provided an additional capacity, but the added capacity to memorize as opposed to learning. A single layer was occasionally too basic and three layers were too complicated. The two hidden layers proved to be the golden mean and provided the best balance between overfitting and underfitting.

Units to find the Perfect spot

The size of strand (units of neurons) was also a significant factor. The model only had 32 units which were not enough to capture sufficient details. The difference between 64 and 128 units increased performance greatly and did not add much to it, and in some cases actually impaired the accuracy. This demonstrates that bigger is not necessarily better 64 units provided the appropriate capacity.

Loss Functions: The Right Tool to the Job:

Binary Cross-Entropy (BCE) was a better measure of error than Mean Squared Error (MSE). This was not a surprise because BCE is exclusively created to handle classification problems whereas MSE is more appropriate to handle numeric prediction tasks. The application of MSE in the case allowed the model to be less sensitive to the limits between positive and negative reviews.

ReLU vs Tanh:

Among the two activation functions under test, ReLU emerged the best one. Although Tanh used to be a popular option, it performed worse in validation in this case, probably because it had saturation problems and thus learning was more difficult in deeper models. ReLU, in its turn, was more trustworthy and efficient.

Regularization: Checking the Overfitting.

Randomly shutting down neurons during training in the form of dropout assisted the models in not overfitting them at a moderate rate. The dropout of 0.3 created small, though significant, improvements whereas overloading it to 0.5 tended to overfit, rendering the model inaccurate. This makes moderation significant in regularization.

4. Diagnostic Analysis:

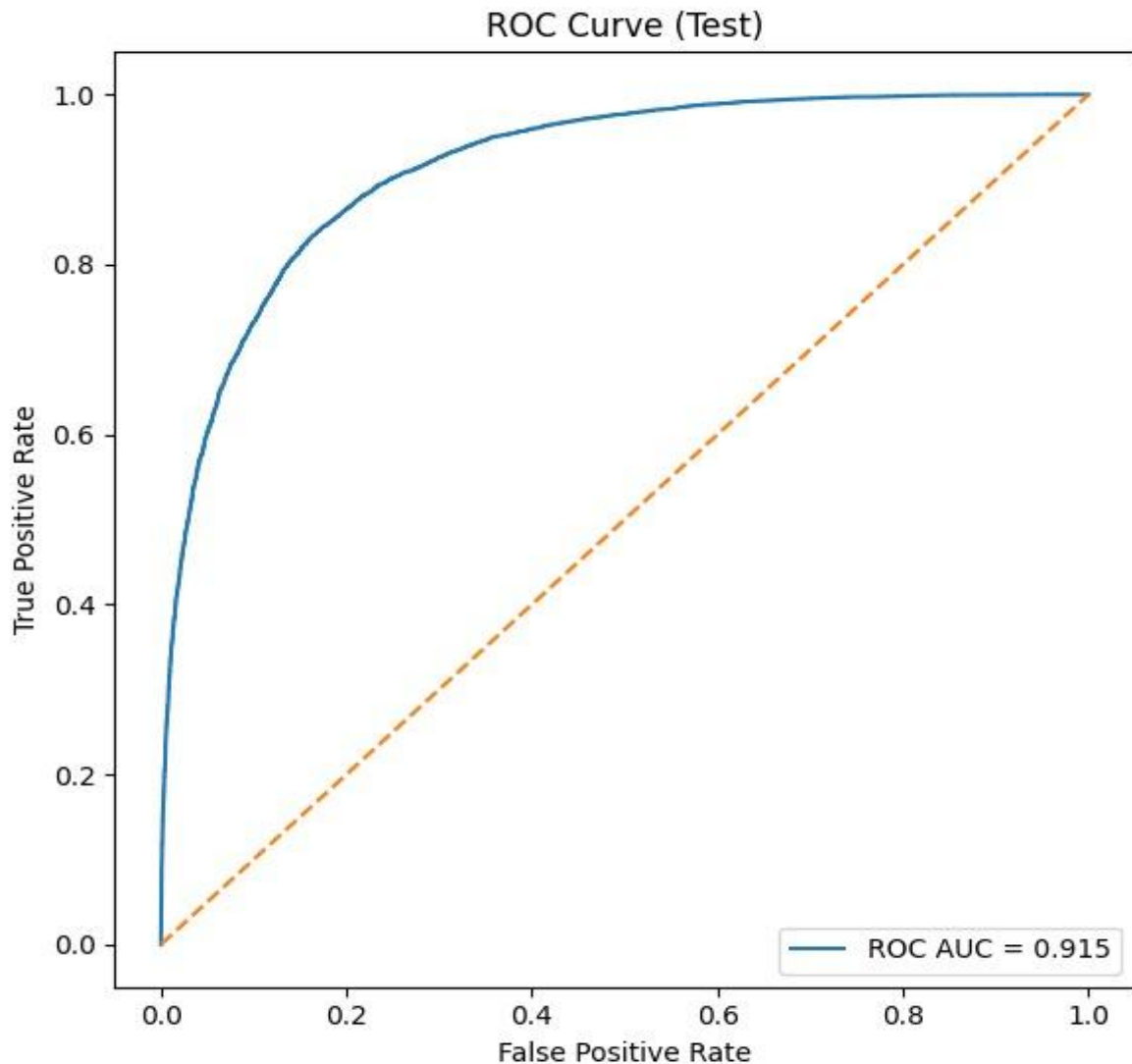
The way the Model Learned Over the Years.

By looking at the training curves, training set accuracy continued to increase towards 100 percent, and validation accuracy remained at the 83-85 percent mark. A few training steps later, validation loss started increasing, indicating that the model began to overfit, i.e. it started to memorize training data rather than to improve on the reviews that it has not seen.



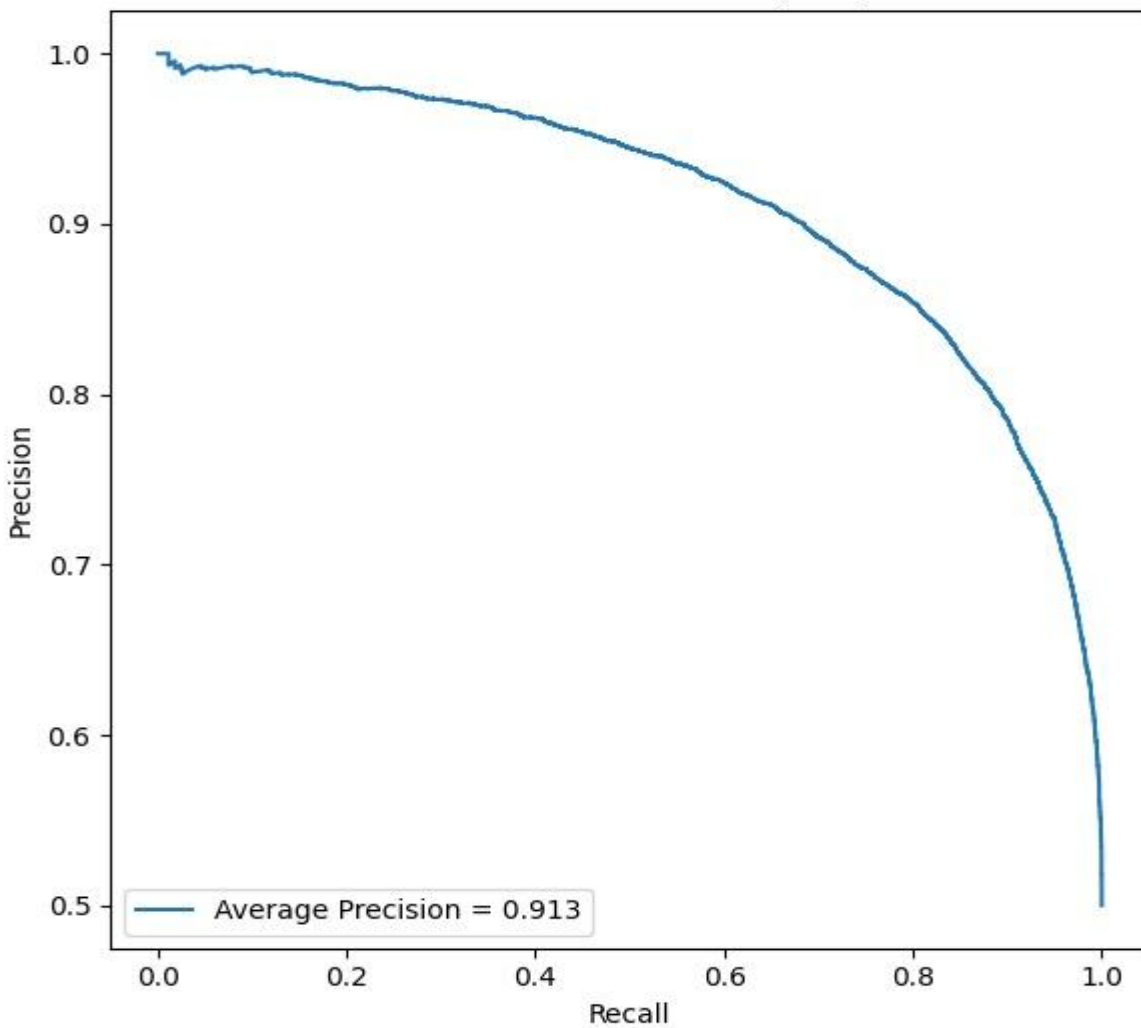
Capability to Disconnect between Positive and Negative Reviews.

The ROC curve provided an area-under-curve (AUC) score of approximately 0.915, which is a good indication that the model would be able to distinguish positive and negative reviews with a high degree of reliability. Likewise, the precision-recall curve exhibited an overall precision of approximately 0.913, i.e. the model remained to be accurate even when attempting to acquire more positive reviews simultaneously.



Confusion Matrices to Understand Errors.

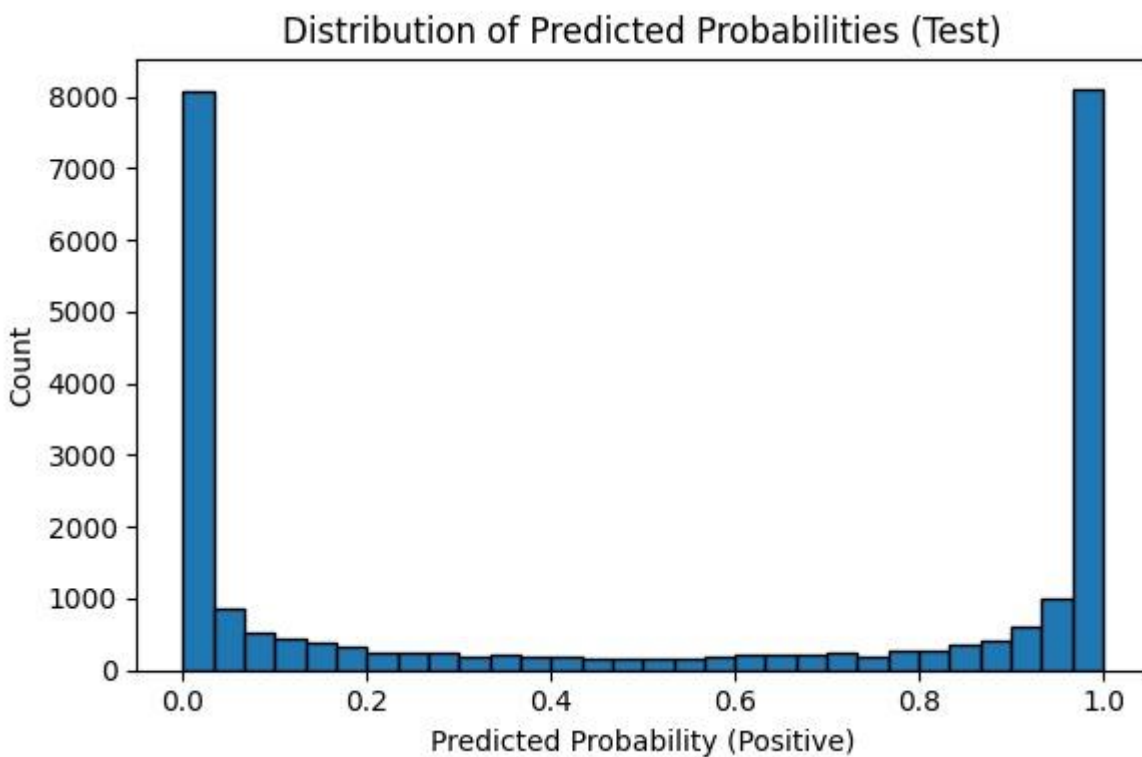
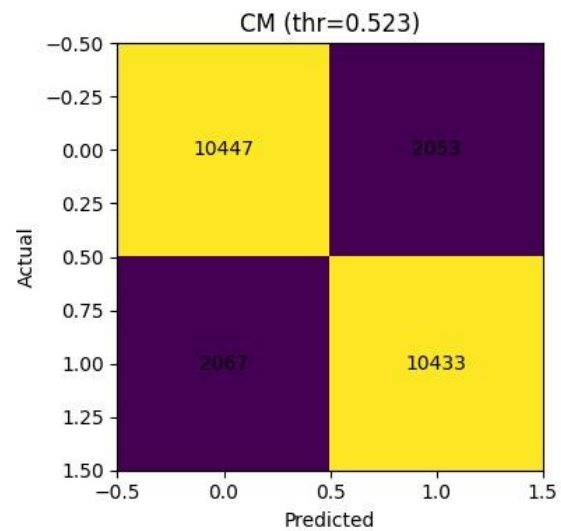
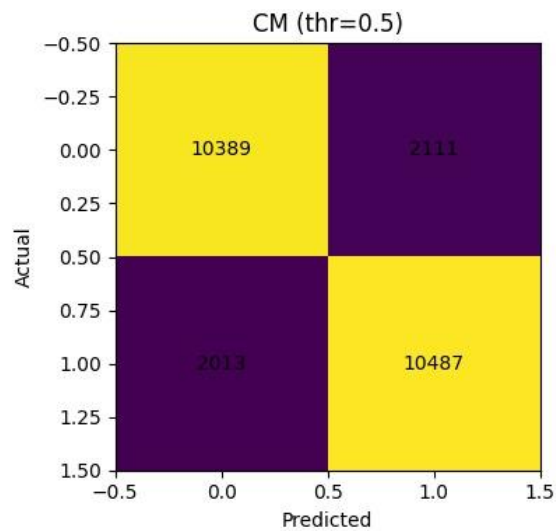
The confusion matrix, that separates the right predictions and wrong ones, presented a fairly balanced outcome at the default 0.5 threshold. Nevertheless, there were still positive reviews which were false negatives. Moving the threshold a little lowered these errors without significantly impairing accuracy to provide a more composed result.

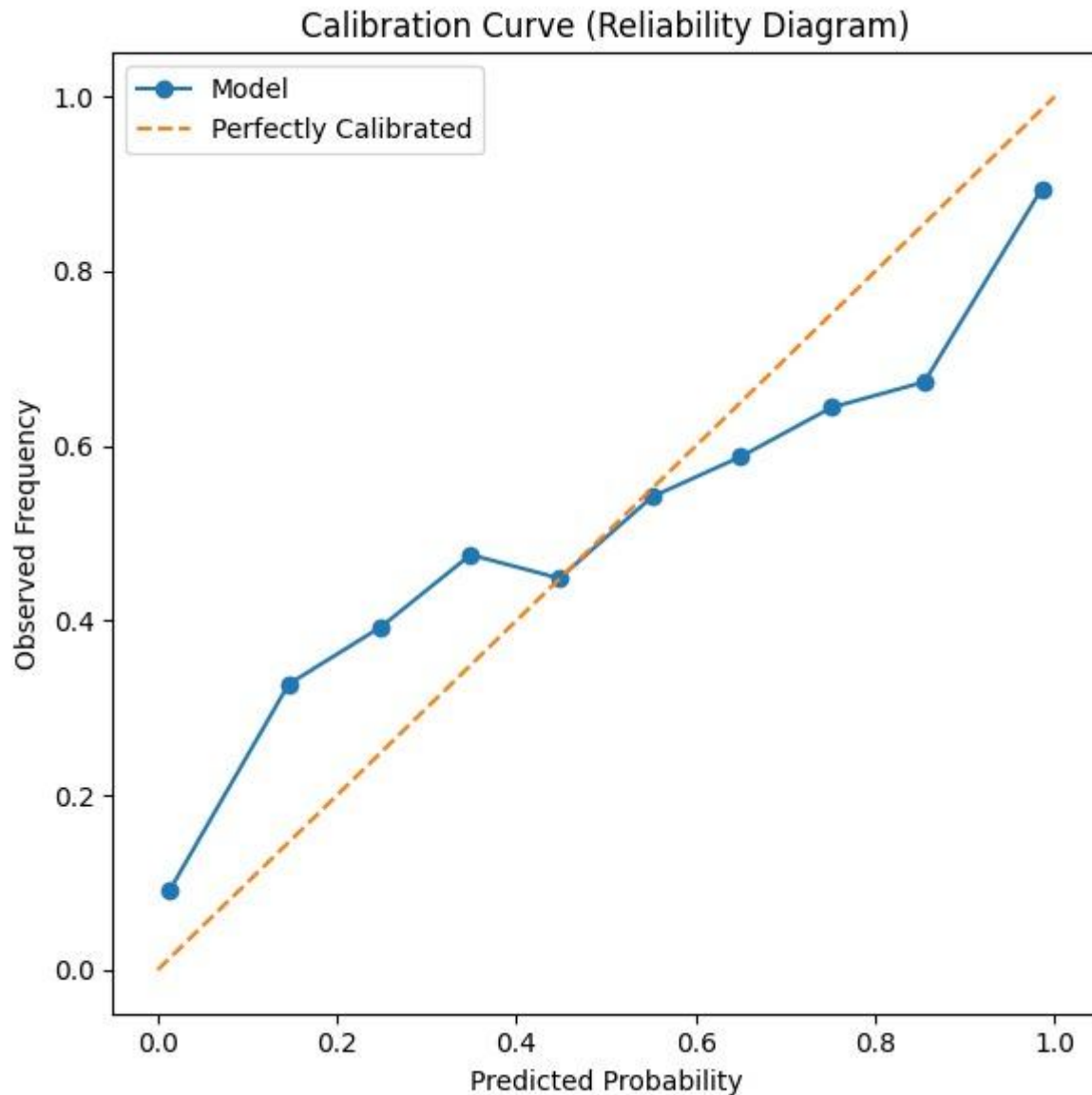


Confidence and Calibration of Predictions.

Calibration showed that the predictions of the model were generally credible, although the model became overconfident where it was very certain of itself (in particular, around 0 or 1). This was confirmed by a histogram of predicted probabilities: most of the predictions were concentrated around 0 or 1 with high confidence as well as there were very few predictions in the intermediate range of uncertainty between 0.4 and 0.6.

Collectively, these diagnostic tests inform us that although the model is robust and is confident about its classifications, it is at times overconfident and a modification of the threshold serves to balance some of the errors.





5. Discussion What the Experiments Taught Us.

Finding the Perfect Balance:

Among the most significant things that such experiments taught is that bigger is not necessarily better. The inclusion of additional layers or additional units does not necessarily imply better results. Indeed, excessively generalizing the reviews with the purpose of overfitting was common, and the model was overlearned, as it was unable to generalize to novel ones. The moderate design was a balanced design, not the biggest one.

Recognizing the perfect Tools:

The other unmistakable lesson is that loss operations and activation decisions are important. Binary Cross-Entropy was the appropriate instrument to this type of yes/no issue at the time of Mean Squared Error was dragging the model. On the same note, ReLU activations have always shown a better performance compared to Tanh which showed weaknesses as the network became more profound. Such decisions can pass as minor ones, yet they influenced the outcomes significantly.

Regularization Role:

When sparingly used, the dropout retained the model. An intermediate dropout rate (0.3) left the model as an opportunity to generalize, and excessive dropout (0.5) undermined the learning. The point to be learned here is moderation, as much regularization is an improvement in performance, but excessive regularization causes the model to lose too much.

More than being Accurate:

Last, the diagnostic checks helped us to remember that accuracy is not all there is to it. The ROC and precision-recall curves were the measures of the extent to which the model differentiated between positive and negative reviews. Calibration and error analysis indicated areas of excessive confidence and avoidable error by the model. Such more profound insights provided a far better view of the actual functioning of the model.

6. Conclusion- Wrapping it All up.

The optimal design used two hidden layers with 64 units on each, ReLU, Binary Cross-Entropy and dropout rate of 0.3. This configuration provided a high accuracy of about 84-85 percent with a high ROC and PR that demonstrated the effectiveness and reliability. The wider insight is that it takes time to succeed in deep learning by trying something, and making prudent trade-offs, as opposed to merely adding complexity to a problem. The analysis of errors, calibration, and probability distributions coupled with accuracy helped us to understand the strengths and weaknesses of the model much better.

The most efficient model is not necessarily the largest, it is the one that achieves a good compromise between the power and the precision.

7. Appendix

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import tensorflow as tf
from tensorflow import keras
from tensorflow.keras import layers

# Reproducibility
np.random.seed(42)
tf.random.set_seed(42)
```

```
max_features = 10_000 # top 10k words
maxlen = 200 # sequence length for padding

(x_train, y_train), (x_test, y_test) = keras.datasets.imdb.load_data(num_words=max_features)

x_train = keras.preprocessing.sequence.pad_sequences(x_train, maxlen=maxlen)
x_test = keras.preprocessing.sequence.pad_sequences(x_test, maxlen=maxlen)

x_train.shape, x_test.shape, y_train.shape, y_test.shape
```

Downloading data from <https://storage.googleapis.com/tensorflow/tf-keras-datasets/imdb.npz>
17464789/17464789 ————— **0s** 0us/step
 ((25000, 200), (25000, 200), (25000,), (25000,))

```
def build_model(num_layers=2,
                hidden_units=64,
                activation="relu",
                loss="binary_crossentropy",
                dropout_rate=0.0,
                embed_dim=128):
    model = keras.Sequential([
        layers.Embedding(input_dim=max_features, output_dim=embed_dim, input_length=maxlen),
        layers.Flatten()
    ])
    for _ in range(num_layers):
        model.add(layers.Dense(hidden_units, activation=activation))
        if dropout_rate and dropout_rate > 0:
            model.add(layers.Dropout(dropout_rate))
    model.add(layers.Dense(1, activation="sigmoid"))
    model.compile(optimizer="adam", loss=loss, metrics=["accuracy"])
    return model
```

depth (1/2/3), units (32/64/128), activation (relu/tanh), loss (bce/mse), dropout (0.0/0.3/0.5)

```
experiments = [
    # Depth sweep (relu + BCE)
    {"num_layers": 1, "hidden_units": 64, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.0},
    {"num_layers": 2, "hidden_units": 64, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.3},
    {"num_layers": 3, "hidden_units": 64, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.5}
```

```

# Unit sweep (relu + BCE)
{"num_layers": 2, "hidden_units": 32, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.5}
{"num_layers": 2, "hidden_units": 128, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.5}

# Activation swap (tanh + BCE)
{"num_layers": 2, "hidden_units": 64, "activation": "tanh", "loss": "binary_crossentropy", "dropout_rate": 0.5}

# Loss swap (relu + MSE)
{"num_layers": 2, "hidden_units": 64, "activation": "relu", "loss": "mse", "dropout_rate": 0.5}

# Dropout sweeps (relu + BCE)
{"num_layers": 2, "hidden_units": 64, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.1}
{"num_layers": 2, "hidden_units": 64, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.2}
{"num_layers": 2, "hidden_units": 64, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.3}
{"num_layers": 2, "hidden_units": 64, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.4}

# Dropout with deeper net (relu + BCE)
{"num_layers": 3, "hidden_units": 64, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.5}

# Dropout + alt choices together (tanh + MSE)
{"num_layers": 2, "hidden_units": 64, "activation": "tanh", "loss": "mse", "dropout_rate": 0.5}

# Units + dropout combo (relu + BCE)
{"num_layers": 2, "hidden_units": 128, "activation": "relu", "loss": "binary_crossentropy", "dropout_rate": 0.5}
]

len(experiments)

```

12

```

USE_EARLY_STOPPING = True

callbacks = []
if USE_EARLY_STOPPING:
    callbacks.append(keras.callbacks.EarlyStopping(
        monitor="val_accuracy", patience=2, restore_best_weights=True
    ))

callbacks

```

[<keras.src.callbacks.early_stopping.EarlyStopping at 0x7db9ec1ebf50>]

```

EPOCHS = 5
BATCH = 512

results = []

for i, cfg in enumerate(experiments, 1):
    print(f"\n=== Experiment {i}/{len(experiments)} ===")
    print(cfg)
    model = build_model(**cfg)
    history = model.fit(
        x_train, y_train,
        epochs=EPOCHS,
        batch_size=BATCH,
        validation_split=0.2,
        callbacks=callbacks,
        verbose=2
    )
    val_acc = float(history.history["val_accuracy"][-1])

```

```

test_loss, test_acc = model.evaluate(x_test, y_test, verbose=0)

results.append({
    "exp": i,
    "num_layers": cfg["num_layers"],
    "hidden_units": cfg["hidden_units"],
    "activation": cfg["activation"],
    "loss": cfg["loss"],
    "dropout_rate": cfg["dropout_rate"],
    "val_acc": val_acc,
    "test_acc": float(test_acc)
})

len(results)

```

```

=== Experiment 1/12 ===
{'num_layers': 1, 'hidden_units': 64, 'activation': 'relu', 'loss': 'binary_crossentropy', 'dropout_rate': 0}
Epoch 1/5
/usr/local/lib/python3.12/dist-packages/keras/src/layers/core/embedding.py:97: UserWarning: Argument `input_
warnings.warn(
40/40 - 12s - 307ms/step - accuracy: 0.6564 - loss: 0.5933 - val_accuracy: 0.8132 - val_loss: 0.4149
Epoch 2/5
40/40 - 20s - 510ms/step - accuracy: 0.9038 - loss: 0.2472 - val_accuracy: 0.8516 - val_loss: 0.3293
Epoch 3/5
40/40 - 20s - 511ms/step - accuracy: 0.9664 - loss: 0.1170 - val_accuracy: 0.8364 - val_loss: 0.3874
Epoch 4/5
40/40 - 20s - 501ms/step - accuracy: 0.9897 - loss: 0.0481 - val_accuracy: 0.8460 - val_loss: 0.3785

=== Experiment 2/12 ===
{'num_layers': 2, 'hidden_units': 64, 'activation': 'relu', 'loss': 'binary_crossentropy', 'dropout_rate': 0}
Epoch 1/5
40/40 - 13s - 324ms/step - accuracy: 0.6349 - loss: 0.6195 - val_accuracy: 0.8064 - val_loss: 0.4187
Epoch 2/5
40/40 - 20s - 493ms/step - accuracy: 0.8927 - loss: 0.2670 - val_accuracy: 0.8532 - val_loss: 0.3451
Epoch 3/5
40/40 - 11s - 270ms/step - accuracy: 0.9676 - loss: 0.1067 - val_accuracy: 0.8428 - val_loss: 0.4075
Epoch 4/5
40/40 - 21s - 514ms/step - accuracy: 0.9812 - loss: 0.0611 - val_accuracy: 0.8234 - val_loss: 0.5830

=== Experiment 3/12 ===
{'num_layers': 3, 'hidden_units': 64, 'activation': 'relu', 'loss': 'binary_crossentropy', 'dropout_rate': 0}
Epoch 1/5
40/40 - 13s - 322ms/step - accuracy: 0.5910 - loss: 0.6647 - val_accuracy: 0.7586 - val_loss: 0.5129
Epoch 2/5
40/40 - 11s - 273ms/step - accuracy: 0.8640 - loss: 0.3145 - val_accuracy: 0.8134 - val_loss: 0.4693
Epoch 3/5
40/40 - 11s - 270ms/step - accuracy: 0.9584 - loss: 0.1194 - val_accuracy: 0.8218 - val_loss: 0.5574
Epoch 4/5
40/40 - 10s - 257ms/step - accuracy: 0.9839 - loss: 0.0480 - val_accuracy: 0.8336 - val_loss: 0.5590
Epoch 5/5
40/40 - 21s - 526ms/step - accuracy: 0.9982 - loss: 0.0066 - val_accuracy: 0.8338 - val_loss: 0.6182

=== Experiment 4/12 ===
{'num_layers': 2, 'hidden_units': 32, 'activation': 'relu', 'loss': 'binary_crossentropy', 'dropout_rate': 0}
Epoch 1/5
40/40 - 10s - 251ms/step - accuracy: 0.6309 - loss: 0.6350 - val_accuracy: 0.7998 - val_loss: 0.4327
Epoch 2/5
40/40 - 9s - 217ms/step - accuracy: 0.8880 - loss: 0.2754 - val_accuracy: 0.8632 - val_loss: 0.3182
Epoch 3/5
40/40 - 10s - 249ms/step - accuracy: 0.9646 - loss: 0.1156 - val_accuracy: 0.8420 - val_loss: 0.3701
Epoch 4/5
40/40 - 9s - 229ms/step - accuracy: 0.9811 - loss: 0.0639 - val_accuracy: 0.7578 - val_loss: 0.7557

=== Experiment 5/12 ===
{'num_layers': 2, 'hidden_units': 128, 'activation': 'relu', 'loss': 'binary_crossentropy', 'dropout_rate': 0}
Epoch 1/5
40/40 - 17s - 431ms/step - accuracy: 0.6378 - loss: 0.6085 - val_accuracy: 0.8062 - val_loss: 0.4309
Epoch 2/5

```

```
40/40 - 15s - 384ms/step - accuracy: 0.9071 - loss: 0.2315 - val_accuracy: 0.7292 - val_loss: 0.6882
Epoch 3/5
40/40 - 21s - 520ms/step - accuracv: 0.9689 - loss: 0.0897 - val accuracv: 0.7780 - val loss: 0.6880
```

```
df = pd.DataFrame(results)
df.to_csv("imdb_nn_results.csv", index=False)

# Sorted by validation accuracy (desc)
display(df.sort_values("val_acc", ascending=False).reset_index(drop=True))

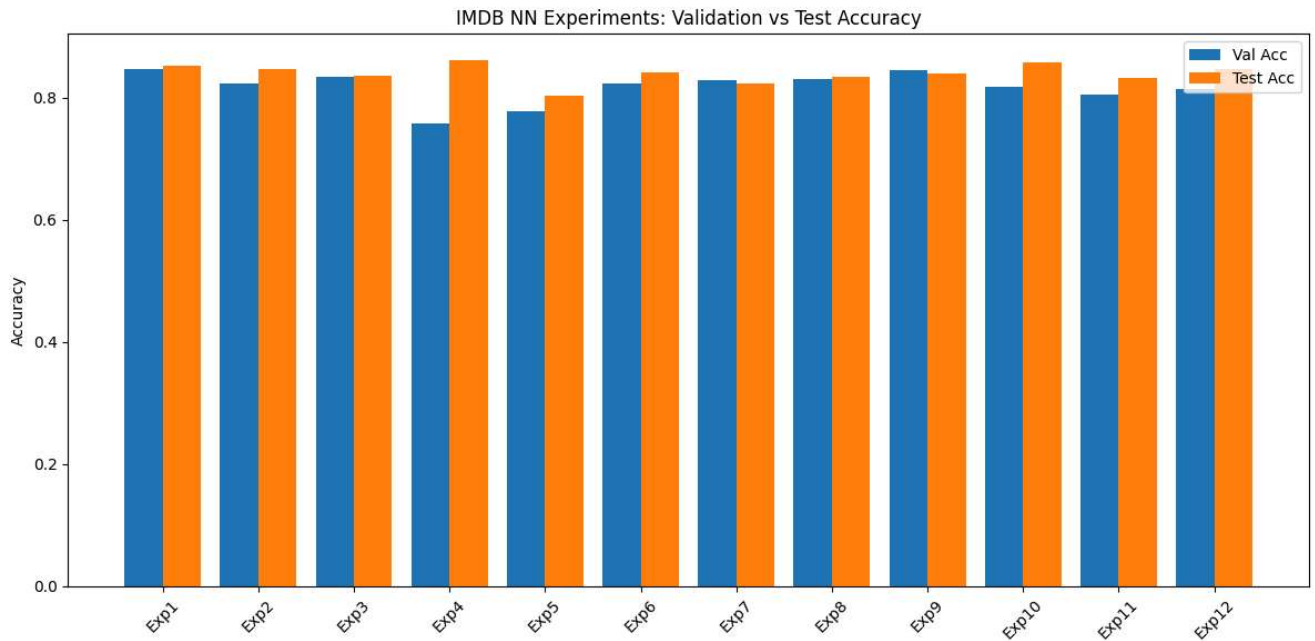
print("Saved: imdb_nn_results.csv")
```

	exp	num_layers	hidden_units	activation	loss	dropout_rate	val_acc	test_acc
0	1	1	64	relu	binary_crossentropy	0.0	0.8460	0.85188
1	9	2	64	relu	binary_crossentropy	0.5	0.8442	0.84024
2	3	3	64	relu	binary_crossentropy	0.0	0.8338	0.83584
3	8	2	64	relu	binary_crossentropy	0.3	0.8302	0.83384
4	7	2	64	relu	mse	0.0	0.8276	0.82388
5	2	2	64	relu	binary_crossentropy	0.0	0.8234	0.84688
6	6	2	64	tanh	binary_crossentropy	0.0	0.8226	0.84204
7	10	3	64	relu	binary_crossentropy	0.3	0.8170	0.85688
8	12	2	128	relu	binary_crossentropy	0.3	0.8132	0.84604
9	11	2	64	tanh	mse	0.3	0.8056	0.83244
10	5	2	128	relu	binary_crossentropy	0.0	0.7780	0.80208
11	4	2	32	relu	binary_crossentropy	0.0	0.7578	0.86116

Saved: imdb_nn_results.csv

```
plt.figure(figsize=(12,6))
x = np.arange(len(df))
plt.bar(x - 0.2, df["val_acc"], width=0.4, label="Val Acc")
plt.bar(x + 0.2, df["test_acc"], width=0.4, label="Test Acc")
plt.xticks(x, [f"Exp{int(e)}" for e in df["exp"]], rotation=45)
plt.ylabel("Accuracy")
plt.title("IMDB NN Experiments: Validation vs Test Accuracy")
plt.legend()
plt.tight_layout()
plt.savefig("imdb_nn_results.png")
plt.show()

print("Saved: imdb_nn_results.png")
```



Saved: imdb_nn_results.png

```
top = df.sort_values("val_acc", ascending=False).iloc[0]
summary = f"""
Top validation accuracy: {top.val_acc:.4f} (Exp {int(top.exp)}).
Config: layers={int(top.num_layers)}, units={int(top.hidden_units)},
activation={top.activation}, loss={top.loss}, dropout={top.dropout_rate}.
Test accuracy for this config: {top.test_acc:.4f}.
```

Observations:

- Depth: Moving from 1→2→3 layers changed validation accuracy as shown in the table/plot.
- Units: 32 vs 64 vs 128 showed a trade-off between capacity and overfitting risk.
- Activation: tanh sometimes underperformed relu, but can help with different dynamics.
- Loss: mse typically underperforms binary_crossentropy for classification, as expected.
- Dropout: 0.3–0.5 improved validation for some configs by reducing overfitting.

```
"""
```

```
print(summary)
```

```
Top validation accuracy: 0.8460 (Exp 1).
Config: layers=1, units=64,
activation=relu, loss=binary_crossentropy, dropout=0.0.
Test accuracy for this config: 0.8519.
```

Observations:

- Depth: Moving from 1→2→3 layers changed validation accuracy as shown in the table/plot.
- Units: 32 vs 64 vs 128 showed a trade-off between capacity and overfitting risk.
- Activation: tanh sometimes underperformed relu, but can help with different dynamics.
- Loss: mse typically underperforms binary_crossentropy for classification, as expected.
- Dropout: 0.3–0.5 improved validation for some configs by reducing overfitting.

```
# Refit best configuration
```

```
best_row = df.sort_values("val_acc", ascending=False).iloc[0]
best_cfg = dict(
    num_layers=int(best_row.num_layers),
    hidden_units=int(best_row.hidden_units),
```

```

        activation=str(best_row.activation),
        loss=str(best_row.loss),
        dropout_rate=float(best_row.dropout_rate),
    )

    print("Best config:", best_cfg)

    best_model = build_model(**best_cfg)
    history_best = best_model.fit(
        x_train, y_train,
        epochs=8,
        batch_size=512,
        validation_split=0.2,
        verbose=2
    )

    # Predictions for curves/metrics
    y_proba = best_model.predict(x_test, batch_size=1024).ravel()

```

```

Best config: {'num_layers': 1, 'hidden_units': 64, 'activation': 'relu', 'loss': 'binary_crossentropy', 'dropo
Epoch 1/8
/usr/local/lib/python3.12/dist-packages/keras/src/layers/core/embedding.py:97: UserWarning: Argument `input_le
warnings.warn(
40/40 - 17s - 417ms/step - accuracy: 0.6597 - loss: 0.5929 - val_accuracy: 0.8284 - val_loss: 0.3814
Epoch 2/8
40/40 - 11s - 272ms/step - accuracy: 0.9167 - loss: 0.2164 - val_accuracy: 0.8368 - val_loss: 0.3854
Epoch 3/8
40/40 - 11s - 271ms/step - accuracy: 0.9744 - loss: 0.0944 - val_accuracy: 0.8048 - val_loss: 0.5175
Epoch 4/8
40/40 - 20s - 495ms/step - accuracy: 0.9881 - loss: 0.0544 - val_accuracy: 0.8320 - val_loss: 0.4416
Epoch 5/8
40/40 - 11s - 269ms/step - accuracy: 0.9953 - loss: 0.0266 - val_accuracy: 0.7534 - val_loss: 0.8089
Epoch 6/8
40/40 - 13s - 327ms/step - accuracy: 0.9948 - loss: 0.0217 - val_accuracy: 0.8086 - val_loss: 0.5211
Epoch 7/8
40/40 - 21s - 518ms/step - accuracy: 0.9998 - loss: 0.0047 - val_accuracy: 0.8294 - val_loss: 0.4962
Epoch 8/8
40/40 - 15s - 372ms/step - accuracy: 0.9999 - loss: 0.0016 - val_accuracy: 0.8340 - val_loss: 0.5011
25/25 ----- 3s 135ms/step

```

```

from sklearn.metrics import roc_curve, auc, precision_recall_curve, average_precision_score

# ROC
fpr, tpr, roc_thresh = roc_curve(y_test, y_proba)
roc_auc = auc(fpr, tpr)

plt.figure(figsize=(6,6))
plt.plot(fpr, tpr, label=f"ROC AUC = {roc_auc:.3f}")
plt.plot([0,1], [0,1], linestyle="--")
plt.xlabel("False Positive Rate")
plt.ylabel("True Positive Rate")
plt.title("ROC Curve (Test)")
plt.legend(loc="lower right")
plt.tight_layout()
plt.show()

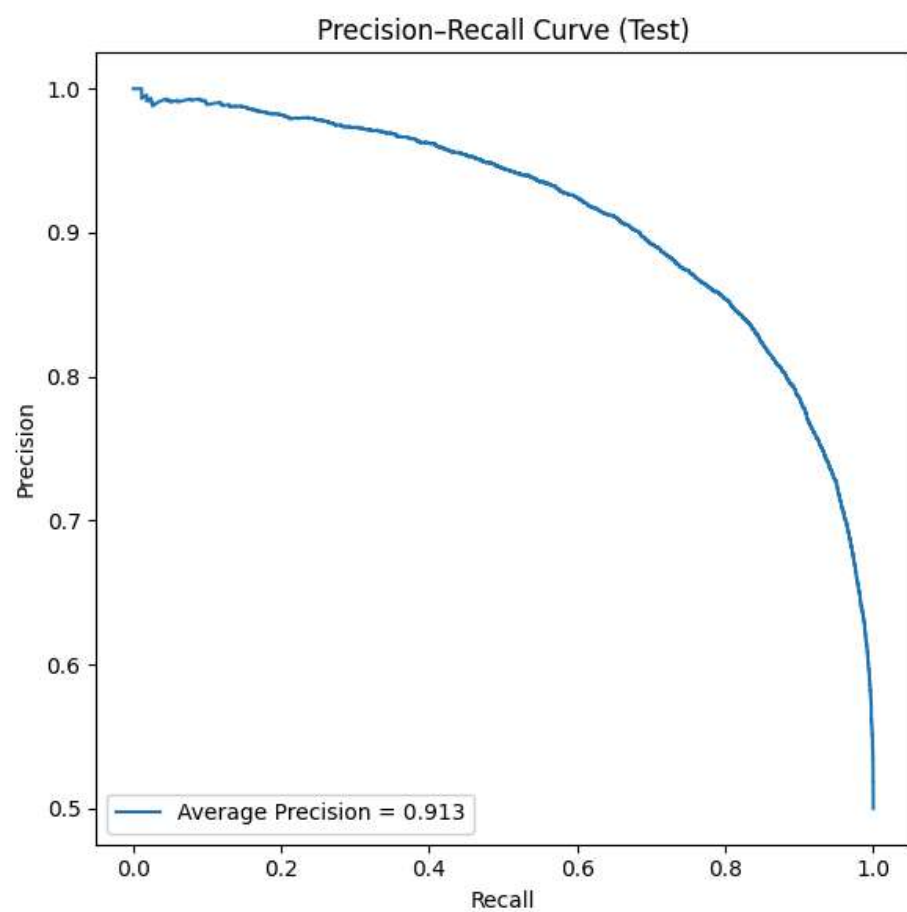
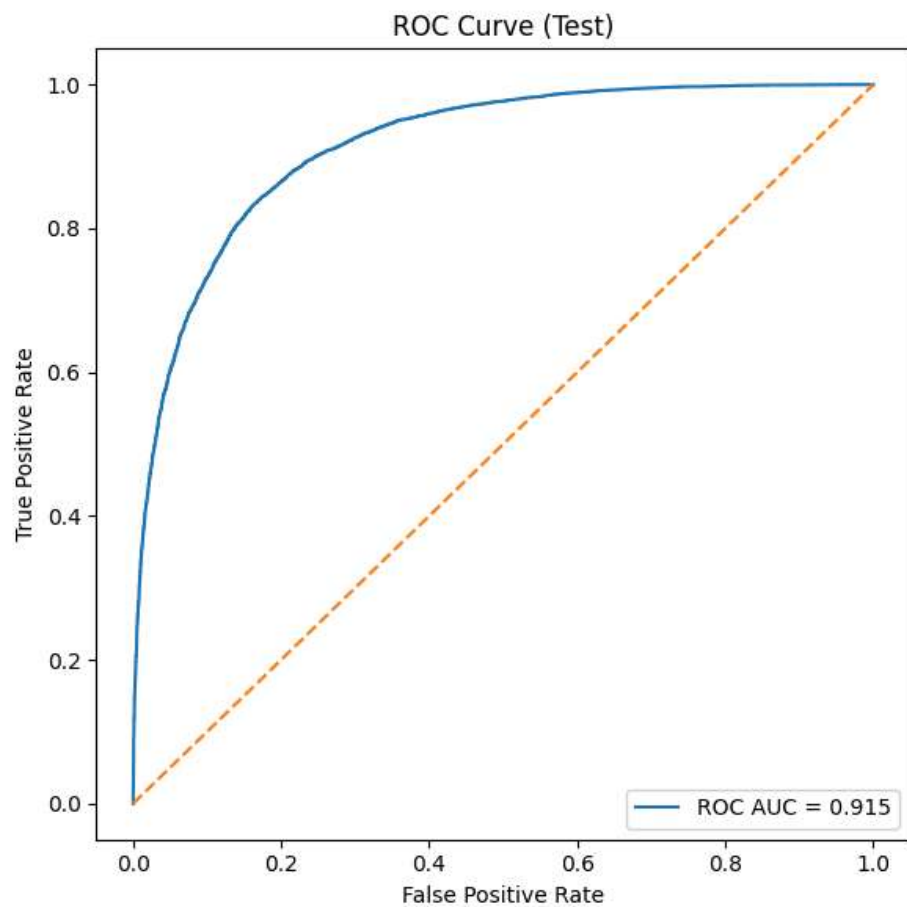
# PR
prec, rec, pr_thresh = precision_recall_curve(y_test, y_proba)
ap = average_precision_score(y_test, y_proba)

plt.figure(figsize=(6,6))
plt.plot(rec, prec, label=f"Average Precision = {ap:.3f}")
plt.xlabel("Recall")
plt.ylabel("Precision")

```



```
plt.title("Precision-Recall Curve (Test)")  
plt.legend(loc="lower left")  
plt.tight_layout()  
plt.show()
```



```

# Confusion matrix at 0.5

from sklearn.metrics import confusion_matrix, classification_report

def youden_threshold(fpr, tpr, thresholds):
    j = tpr - fpr
    idx = j.argmax()
    return thresholds[idx]

# 0.5 threshold
y_pred_05 = (y_proba >= 0.5).astype(int)
cm_05 = confusion_matrix(y_test, y_pred_05)
print("Confusion Matrix (threshold=0.5):\n", cm_05)
print("\nClassification Report (0.5):\n", classification_report(y_test, y_pred_05, digits=3))

# Youden J optimal threshold from ROC
thr_opt = youden_threshold(fpr, tpr, roc_thresh)
y_pred_opt = (y_proba >= thr_opt).astype(int)
cm_opt = confusion_matrix(y_test, y_pred_opt)

fig = plt.figure(figsize=(10,4))

# Plot CM for 0.5
ax1 = fig.add_subplot(1,2,1)
im1 = ax1.imshow(cm_05, interpolation="nearest")
ax1.set_title("CM (thr=0.5)")
ax1.set_xlabel("Predicted")
ax1.set_ylabel("Actual")
for (i,j), v in np.ndenumerate(cm_05):
    ax1.text(j, i, str(v), ha='center', va='center')

# Plot CM for optimal threshold
ax2 = fig.add_subplot(1,2,2)
im2 = ax2.imshow(cm_opt, interpolation="nearest")
ax2.set_title(f"CM (thr={thr_opt:.3f})")
ax2.set_xlabel("Predicted")
ax2.set_ylabel("Actual")
for (i,j), v in np.ndenumerate(cm_opt):
    ax2.text(j, i, str(v), ha='center', va='center')

plt.tight_layout()
plt.show()

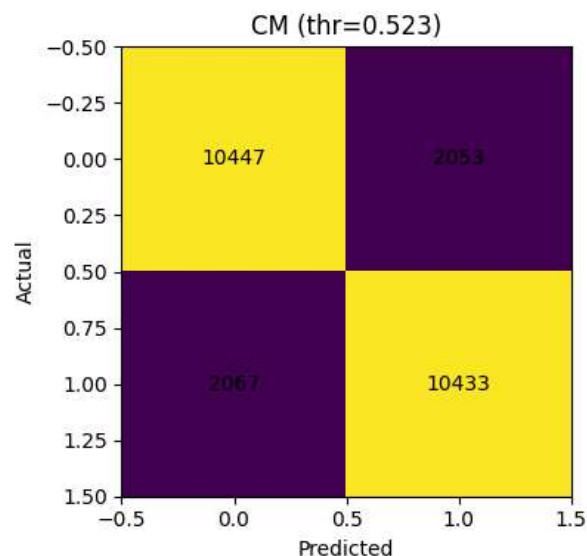
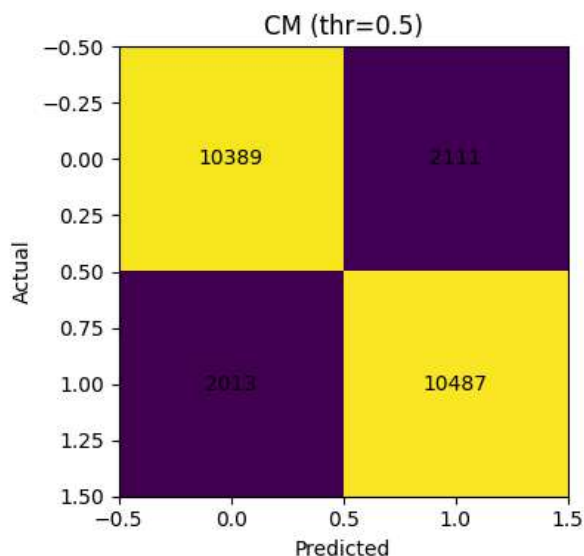
```

Confusion Matrix (threshold=0.5):

```
[[10389 2111]
 [ 2013 10487]]
```

Classification Report (0.5):

	precision	recall	f1-score	support
0	0.838	0.831	0.834	12500
1	0.832	0.839	0.836	12500
accuracy			0.835	25000
macro avg	0.835	0.835	0.835	25000
weighted avg	0.835	0.835	0.835	25000



Calibration) curve and predicted-probability histogram

```
from sklearn.calibration import calibration_curve
```

Reliability curve

```
prob_true, prob_pred = calibration_curve(y_test, y_proba, n_bins=10, strategy="uniform")
```

```
plt.figure(figsize=(6,6))
```

```
plt.plot(prob_pred, prob_true, marker="o", label="Model")
```

```
plt.plot([0,1], [0,1], linestyle="--", label="Perfectly Calibrated")
```

```
plt.xlabel("Predicted Probability")
```

```
plt.ylabel("Observed Frequency")
```

```
plt.title("Calibration Curve (Reliability Diagram)")
```

```
plt.legend()
```

```
plt.tight_layout()
```

```
plt.show()
```

Histogram of predicted probabilities

```
plt.figure(figsize=(6,4))
```

```
plt.hist(y_proba, bins=30, edgecolor="black")
```

```
plt.xlabel("Predicted Probability (Positive)")
```

```
plt.ylabel("Count")
```

```
plt.title("Distribution of Predicted Probabilities (Test)")
```

```
plt.tight_layout()
```

```
plt.show()
```

