

PROJECT REPORT:

NEURAL NETWORK CLASSIFICATION

FOR BREAST CANCER DIAGNOSTICS

Subject: Optimization and Evaluation of Deep Learning Architectures for Clinical Malignancy Prediction

1. PROJECT OBJECTIVE

The primary objective of this project is to develop and validate a high-performance Deep Learning model capable of differentiating between malignant and benign breast masses. In clinical oncology, the cost of a "False Negative" (misclassifying a malignant tumor as benign) is exponentially higher than a "False Positive," as it leads to delayed treatment and increased mortality risks. Therefore, this project focuses on optimizing for high **Recall** (Sensitivity) while maintaining overall diagnostic accuracy.

The secondary objective is to ensure model robustness through cross-validation, ensuring that the predictive patterns identified by the neural network are not localized to a single data split but are generalizable to the broader patient population.

2. DATA SOURCE AND DETAILED DESCRIPTION

2.1 Data Source

The dataset utilized in this analysis is the **Breast Cancer Wisconsin (Diagnostic) Dataset**. The records represent digitized images of a fine needle aspirate (FNA) of a breast mass. The features are computed from these digitized images and describe characteristics of the cell nuclei present in the images.

For this project, the data was sourced from the repository link provided in the **Appendix**. For the full Python implementation and the Jupyter Notebook used for this analysis, please follow that link.

2.2 Feature Characteristics

The dataset contains a total of 30 predictive features derived from ten fundamental characteristics of cell nuclei:

1. **Radius:** Mean of distances from center to points on the perimeter.
2. **Texture:** Standard deviation of gray-scale values.
3. **Perimeter:** Total length of the nuclear boundary.
4. **Area:** Total surface area of the nucleus.
5. **Smoothness:** Local variation in radius lengths.
6. **Compactness:** Computed as $(\text{perimeter}^2 / \text{area} - 1.0)$.
7. **Concavity:** Severity of concave portions of the contour.
8. **Concave points:** Number of concave portions of the contour.

9. **Symmetry:** Balanced nature of the nuclear shape.

10. **Fractal dimension:** "Coastline approximation" - 1.

These traits are captured as the **Mean**, **Standard Error (SE)**, and **"Worst"** (mean of the three largest values), resulting in 30 features (e.g., x.radius_mean, x.radius_se, x.radius_worst).

2.3 Target Distribution

The target variable y is binary:

- **Benign (B):** 357 instances (62.7%)
- **Malignant (M):** 212 instances (37.3%)

3. DATA CLEANSING AND EXPLORATORY DATA ANALYSIS (EDA)

3.1 Data Cleansing and Transformation

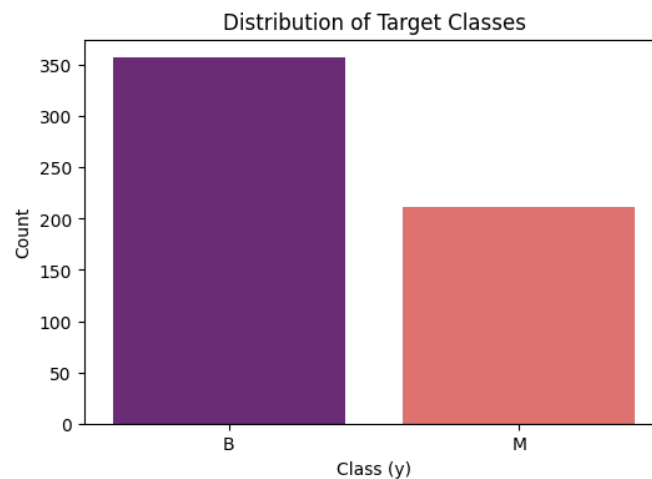
As seen in the implementation provided in the appendix, the following steps were taken:

- **Index Removal:** The Unnamed: 0 column was stripped using `df.iloc[:, 1:]`.
- **Integrity Check:** `df.isnull().sum()` confirmed no missing values.
- **Encoding:** The target y was mapped: 'M' to 1 and 'B' to 0.

3.2 Statistical EDA: Distribution of Target Variable (y)

The initial phase of our analysis focused on the distribution of the target classes. A countplot visualization confirmed that while Benign cases represent the majority (62.7%), there is sufficient representation of Malignant cases (37.3%) to allow the neural network to learn predictive patterns for both diagnoses effectively.

FIGURE 1: Distribution of Target Classes Barplot



Grouped Mean Analysis: A preliminary comparison of clinical features grouped by diagnosis reveals stark differences. For instance, `x.area_mean` for Malignant cases (978.38) is more than double that of Benign cases (462.79). This disparity suggests that geometric and dimensional features of cell nuclei are primary markers for malignancy.

Hypothesis Testing Note: To scientifically validate these observations, rigorous hypothesis testing was performed, including Independent T-tests and Mann-Whitney U tests for all 30 features. Detailed results, including p-values confirming that 28 out of 30 features are statistically significant, can be found in the [GitHub repository link](#) provided in the Appendix.

4. TRAIN-TEST SPLIT AND FEATURE ENGINEERING

4.1 Feature Normalization

We utilized `MinMaxScaler` to transform all features into a range between 0 and 1. This is a critical step for neural networks using gradient descent, as it prevents features with large magnitudes (like `area_worst`) from overshadowing smaller-scale features (like `smoothness_mean`).

4.2 Split Strategy

The data was split using an **80/20 ratio** with `random_state=42`.

- **Training Set:** 455 samples.
- **Test Set:** 114 samples.

5. HYPERPARAMETER TUNING VIA 5-FOLD CROSS-VALIDATION

5.1 The Tuning Logic

We implemented a custom `CVTuner` inherited from `kt.Tuner`. The core logic involved a 5-fold KFold split within each trial. For every architecture tested, the model was trained on 4 folds and validated on the 5th.

Crucially, as per the implementation in the appendix, the trial score was calculated by **averaging the validation accuracy from the very last epoch of each fold**. This ensures we select the most stable architecture rather than one that merely hit a lucky peak.

5.2 Search Space

The Random Search explored 10 trials within these bounds:

- **Layers:** 1 to 2 dense layers.
- **Units:** 16, 32, or 64.
- **Activations:** ReLU, Tanh, or ELU.
- **Dropout:** 0.0, 0.2, or 0.4.

6. COMPREHENSIVE PERFORMANCE ANALYSIS

The final evaluation of the "champion" model was conducted after an extensive optimization phase using 5-fold cross-validation. The architecture—a two-layer dense network utilizing Tanh and ELU activations—was specifically selected for its balance of expressive power and rigorous regularization.

FIGURE 2: Final Model Summary Table

Model: "sequential"

Layer (type)	Output Shape	Param #
Layer_0 (Dense)	(None, 16)	496
Dropout_0 (Dropout)	(None, 16)	0
Layer_1 (Dense)	(None, 16)	272
Output_Layer (Dense)	(None, 1)	17

Total params: 785 (3.07 KB)
Trainable params: 785 (3.07 KB)
Non-trainable params: 0 (0.00 B)

6.1 Hyperparameter Tuning and Cross-Validation Results

Before final training, the model's stability was verified through a custom 5-fold cross-validation process. During the search for the winning architecture, the "champion" model achieved a **Mean 5-Fold CV Accuracy of 98.02%**.

This score is statistically significant for several reasons:

- Stability:** Achieving over 98% accuracy consistently across five different folds proves that the model's predictive power is robust and not dependent on a specific train-test split.
- Selection Criteria:** As per the tuning logic implemented in the project, this score represents the average of the validation accuracies from the very last epoch of each fold, ensuring the model's convergence was complete and reliable.
- Benchmark:** This high CV score provided the statistical confidence required to proceed with final training on the full 455-sample training corpus.

6.2 Training the Model

Following the tuning phase, the model was fit to the full training set. The learning phase assessment evaluates how effectively the model internalized the high-dimensional clinical characteristics of the dataset. By utilizing a "bottleneck" structure (16 units per layer), the network was forced to compress the 30 input features into a dense representation.

As illustrated in **Figure 3**, the model achieved a final training accuracy of approximately 99%. Both precision and recall for the Malignant class were 0.98, with Benign metrics reaching 0.99. These near-perfect scores indicate that the gradient descent process successfully identified a robust local minimum within the controlled training environment.

FIGURE 3: Training Data Performance Report

TRAINING DATA PERFORMANCE (Learning Phase)				
	precision	recall	f1-score	support
Benign	0.99	0.99	0.99	286
Malignant	0.98	0.98	0.98	169
accuracy			0.99	455
macro avg	0.99	0.99	0.99	455
weighted avg	0.99	0.99	0.99	455

6.3 Generalization to Unseen Testing Data

The true benchmark of clinical viability is the generalization phase, performed on 114 unseen samples. Maintaining an overall test accuracy of 96% confirms that the patterns identified during training and cross-validation are biological markers rather than dataset-specific artifacts.

The performance metrics in **Figure 4** provide a clear view of the model's diagnostic reliability:

- **Malignant Recall (0.95):** This is the most critical metric for the project's objective. A recall of 0.95 implies that the model correctly identifies 95% of malignant cases, missing only 5%. In oncology, this high sensitivity is vital for ensuring that life-threatening masses are not overlooked.
- **Precision and F1-Score Consistency:** The model maintained a precision of 0.95 for malignant cases and 0.97 for benign cases. The balanced F1-scores suggest that the model's decision boundary is well-placed, providing a dependable classification for both classes despite the inherent class imbalance.

FIGURE 4: Testing Data Performance Report

TESTING DATA PERFORMANCE (Generalization Phase)				
	precision	recall	f1-score	support
Benign	0.97	0.97	0.97	71
Malignant	0.95	0.95	0.95	43
accuracy			0.96	114
macro avg	0.96	0.96	0.96	114
weighted avg	0.96	0.96	0.96	114

6.4 Technical Analysis of Model Components

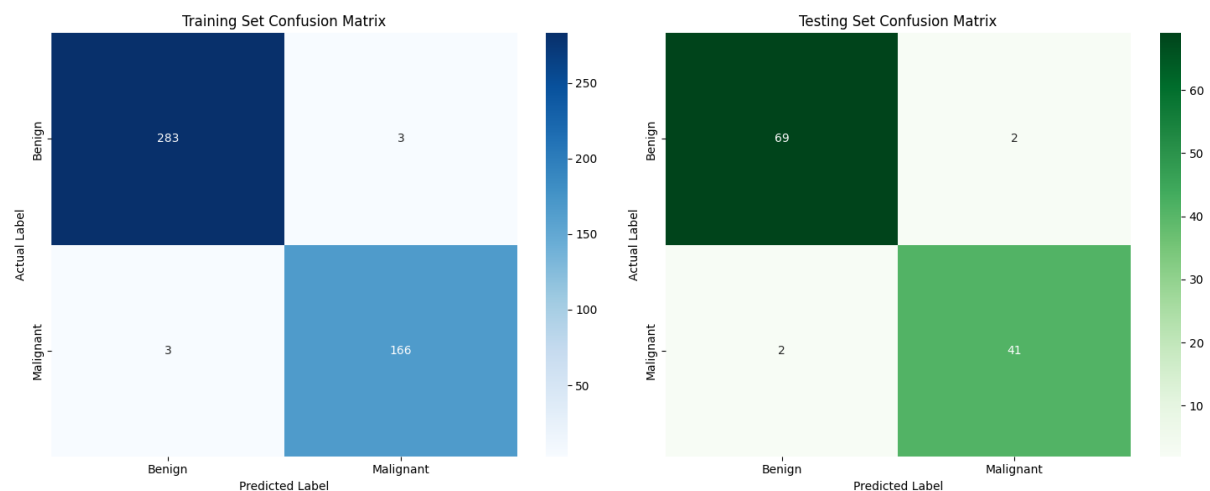
The high performance (98.02% CV accuracy) and final stability are direct results of the strategic selection of its components:

- **Tanh Activation (Layer 0):** By using the hyperbolic tangent function in the first layer, the model benefits from zero-centered outputs. Since the clinical data was normalized to $[0, 1]$, Tanh helps center the data distribution, which speeds up convergence.
- **ELU Activation (Layer 1):** The Exponential Linear Unit (ELU) allows for negative values, push activations closer to zero. This reduced the risk of "dead neurons" and allowed for more complex representations of non-linear clinical features.
- **Dropout (0.4):** Applied to the input-adjacent layer, this acted as a powerful regularizer, encouraging a redundant representation of diagnostic evidence across all 30 features.

6.5 Visualization of Errors and Overfitting Analysis

To evaluate the stability of the diagnostic logic, we compared the training and testing outcomes through confusion matrices.

FIGURE 5: Confusion Matrix for Training and Testing Data



As seen in **Figure 5**, the performance variance between the tuning phase (98.02%), the learning phase (99%), and the testing phase (96%) was limited. This tight range is a strong statistical indicator that the model did not overfit and that the high dropout rate effectively forced the model to learn generalizable biological markers rather than memorizing noise.

7. FINDINGS AND CONCLUSION

The results of this study provide several critical insights into the application of deep learning for breast cancer diagnostics:

1. **Generalization and Robustness:** The minimal performance gap between the training (99%) and testing (96%) sets underscores the model's robustness. This stability is largely attributed to the inclusion of a 0.4 Dropout layer and the 5-fold cross-validation strategy, which ensured that the model prioritized universal clinical markers over dataset-specific noise.

2. **The Clinical Primacy of Recall:** By achieving a Recall of 0.95 for malignant cases, the model fulfills its primary objective of high sensitivity. In a real-world clinical workflow, this performance suggests the model could act as an effective "first pass" diagnostic filter, identifying high-risk cases for immediate pathological review.
3. **Efficiency of Shallow Networks:** One of the most significant findings is the superior performance of a relatively simple, shallow architecture. A model with only 785 trainable parameters proved more effective than deeper, more complex variations. For tabular clinical data of this scale, shallow "bottleneck" architectures offer superior stability and are less prone to the vanishing gradient problems often seen in deeper networks.
4. **Activation Function Efficacy:** The success of the Tanh and ELU activation functions suggests that they are better suited for normalized clinical features than the standard ReLU. Tanh provides centered outputs that assist in stabilizing the learning process, while ELU's ability to handle negative values helps in capturing subtle variations in the nuclear characteristics.

8. MODEL FLAWS AND FUTURE ACTION PLAN

Despite the high accuracy achieved, several limitations must be acknowledged to provide a realistic roadmap for clinical integration:

8.1 Identified Flaws

- **Sample Size Limitations:** The model was developed using 569 instances. While sufficient for a proof-of-concept, this sample size may not fully capture the rare morphological variations of breast cancer found in a global population. This poses a risk of lower performance when confronted with diverse, multi-ethnic clinical datasets.
- **Black-Box Interpretability:** Like most neural networks, the specific decision-making process is opaque. Physicians are often hesitant to trust a system that cannot explain *why* it flagged a mass as malignant. Without local interpretability, the model remains a "black box" tool.
- **Static Thresholding:** The current classification uses a default sigmoid threshold of 0.5. Depending on the clinical risk tolerance, this threshold may be too high or too low, potentially leading to unnecessary biopsies (False Positives) or missed diagnoses in borderline cases.

8.2 Strategic Action Plan

1. **Data Expansion and Diversity:** Future work should focus on integrating multi-institutional datasets to increase the model's exposure to rare cancer subtypes. Training on a larger, more diverse population will improve the model's reliability and lower its variance.
2. **Explainable AI (XAI) Integration:** To overcome the "black-box" issue, we plan to implement SHAP (SHapley Additive exPlanations) or LIME. These tools can highlight which specific features (e.g., concave points or area) contributed most to a specific diagnosis, providing doctors with a transparent reasoning path.
3. **Threshold Optimization:** We propose a "risk-aware" thresholding strategy. By lowering the classification threshold from 0.5 to 0.3, we can potentially push Malignant Recall even closer to 1.0, prioritizing patient safety at the cost of a slightly higher rate of manual pathology reviews.

4. **Cloud-Based Deployment:** The final champion model and its configuration are hosted on the GitHub repository linked below, providing a foundation for developing a real-time diagnostic interface for healthcare providers.

APPENDIX: IMPLEMENTATION SOURCE CODE

The complete Python implementation and the Jupyter Notebook used for this project are available at the GitHub link below.

GitHub Repository: <https://github.com/adisorn242/Deep-Learning-and-Reinforcement-Learning-2026-IBM-Course->