

The Breast Cancer Wisconsin Diagnostic Dataset is a well-known dataset in the field of machine learning and is frequently used for binary classification tasks. The dataset contains 570 samples of breast tissue, where each sample has 30 features describing characteristics of the cell nuclei present in the image of the tissue. The goal is to classify each sample as either benign or malignant based on the features. One popular machine-learning model for this dataset is the Support Vector Machine (SVM). SVMs are powerful classifiers that work by finding the optimal hyperplane that separates the two classes. The hyperplane is chosen to maximize the margin between the two classes, which can improve the model's generalization performance. A project idea using the Breast Cancer Wisconsin Diagnostic Dataset and an SVM could be to develop a diagnostic tool for breast cancer. The SVM could be trained on a subset of the data and then used to classify new samples as either benign or malignant. The tool could be integrated into existing medical software or deployed as a standalone application. The tool could also be used to identify the most important features for diagnosis, which could help in developing new diagnostic tests or treatments for breast cancer. Additionally, the tool could be used to identify false negatives or false positives in existing diagnostic tests, which could help improve their accuracy. Overall, an SVM-based breast cancer diagnostic tool has the potential to improve patient outcomes and advance the field of breast cancer research.

### **Dataset Description**

The Breast Cancer Wisconsin Diagnostic Dataset is a widely used dataset in machine learning and is used for binary classification tasks. The dataset contains 570 samples of breast tissue, where each sample has 30 features describing characteristics of the cell nuclei present in the image of the tissue. The dataset is available through the UCI Machine Learning Repository. The samples are labeled as either "benign" or "malignant" based on the characteristics of the cell nuclei in the image. The features include mean, standard error, and worst (i.e., the largest or most extreme) values of the following characteristics: Radius, Texture, Perimeter, Area, Smoothness, Compactness, Concavity, Concave points, Symmetry, and Fractal dimension. Each feature is represented as a floating-point value. The mean, standard error, and worst values of each feature are provided, resulting in a total of 30 features per sample. The dataset is commonly used for developing machine learning models to predict the likelihood of breast cancer based on the features of the cell nuclei. The goal is to develop models that can accurately classify new tissue samples as either benign or malignant based on their features, which could potentially be used to improve early detection and diagnosis of breast cancer.

## References

- Salama, G. I., Abdelhalim, M., & Zeid, M. A. E. (2012). Breast cancer diagnosis on three different datasets using multi-classifiers. *Breast Cancer (WDBC)*, 32(569), 2.
- Title: "Performance comparison of different classification algorithms for breast cancer dataset." Authors: Ismail, W. S., & Jan, B. Published in: 2014 International Conference on Computer, Communications, and Control Technology, 2014. Link: [Performance comparison of different classification algorithms for breast cancer dataset](#)
- Alshayegi, M. H., Ellethy, H., & Gupta, R. (2022). Computer-aided detection of breast cancer on the Wisconsin dataset: An artificial neural networks approach. *Biomedical Signal Processing and Control*, 71, 103141. [goog](#)

## Comparison

- 69.23% using multiclassifier
- 85% using KNN
- The shallow ANN model showed promising performance with an average accuracy of 90.85%, specificity of 90.72%, sensitivity of 100%, precision of 90.69%, and F1 score of 90.84%.