

About the Breast Cancer Wisconsin (Diagnostic) Dataset


Overview

The **Breast Cancer Wisconsin (Diagnostic) Dataset** is a widely used dataset for binary classification of breast tumors into **benign (B)** or **malignant (M)** categories based on diagnostic measurements. It serves as a benchmark for **machine learning models** in cancer detection and medical research.






Source

This dataset was provided by **Dr. William H. Wolberg** from the University of Wisconsin-Madison and is publicly available on multiple platforms:

 Kaggle Link: [Breast Cancer Wisconsin \(Diagnostic\) Dataset](#)

 UCI Machine Learning Repository: [Breast Cancer Wisconsin \(Diagnostic\)](#)







Dataset Description


-  **Number of Instances:** 569
-  **Number of Features:** 30 (excluding the ID column and target variable)
-  **Target Variable:**
 - **Malignant (M):** 212 instances
 - **Benign (B):** 357 instances
-  **Feature Types:** Real-valued, continuous numerical attributes
-  **Missing Values:** None

Features and Their Role in Classification

The dataset consists of 30 numerical features computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. These features describe characteristics of the cell nuclei present in the image and are categorized into three groups:

✅ **Mean Values:** Describe the average characteristics of the tumor, such as:

-  **Radius (Mean):** Average distance from the center to the perimeter (larger values indicate potential malignancy).
-  **Texture (Mean):** Standard deviation of gray-scale intensity (higher variation may suggest malignancy).
-  **Perimeter & Area (Mean):** Larger values can indicate more aggressive tumors.
-  **Smoothness & Compactness (Mean):** Measure of uniformity and cell cohesion (irregular structures may suggest malignancy).
-  **Concavity & Concave Points (Mean):** Measures of the severity and number of concave portions in tumor contour (more concavity often suggests malignancy).
-  **Symmetry & Fractal Dimension:** Indicators of tumor shape irregularities.


 **Standard Error Values:** Indicate variability in measurements, helping to capture inconsistencies in tumor shape and structure.


⚠️ **Worst Values:** Represent the maximum recorded values for each feature, highlighting extreme cases of tumor growth and irregularity.

These features collectively help in classifying tumors as **benign or malignant** by identifying patterns associated with cancerous growth.

Reference

This dataset has been referenced in:

 K. P. Bennett and O. L. Mangasarian, "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34.


 W. H. Wolberg, W. N. Street, and O. L. Mangasarian. "Breast Cancer Wisconsin (Diagnostic) Data Set."

License

This dataset is licensed under CC BY-NC-SA 4.0.

This dataset is **publicly available** for **educational and research** purposes.

Additional Resources

 **UW CS FTP Server:** To access older versions of the dataset, use:

ftp ftp.cs.wisc.edu

cd math-prog/cpo-dataset/machine-learn/WDBC/