

Bare Demo of IEEEtran.cls for IEEE Computer Society Journals

Michael Shell, *Member, IEEE*, John Doe, *Fellow, OSA*, and Jane Doe, *Life Fellow, IEEE*

Abstract—The abstract goes here.

Index Terms—Computer Society, IEEE, IEEEtran, journal, L^AT_EX, paper, template.

1 INTRODUCTION

2 METHODOLOGY

As the amount of generated per day grows at an exponential rate, brand new technologies have to be developed to cope up with the copious exabytes of data. Machine learning tools provide us with the capabilities to handle both structured and unstructured datasets. These tools can be configured to analyze patterns inherent in the data and make accurate predictions based on the information obtained. This concept is a reality for almost all sectors today. As per a 2020 Stanford study, the amount of healthcare data generated will be around 2,314 exabytes with a steady growth of 48%. The pipeline developed for this project has been depicted in Figure 1. The remainder of the section describes each individual step in detail.

2.1 Obtaining Raw data

In order to obtain distinct yet comparable subjects, a cohort dataset of 89 patients was selected in this study. The dataset consisted of four intrinsic molecular subtypes of breast cancer which are contrasted on the genes a cancerous cell expresses. The dataset has been described in Table 1.

For each of the patient, a CT scan was conducted to obtain cross-sectional images of the hypothesised tumor location. CT scans provide a more detailed description of

the patients condition by increasing the radiation level the patient is exposed to. Once the scan is completed three views are obtained namely, Axial, Sagittal and Coronal. A sample Axial view has been displayed in Figure 2 with the distinct grey circular mass on the right depicting the tumor. DICOM (Digital Imaging and Communication in Medicine) images were obtained after the scan. For each patient 323 new studies were conducted with each study have 384 series which corresponded to 466 instances or images of the scan. Even though DICOM files are a standard format for medical imaging, NRRD (Nearly Raw Raster Data) files are anonymized and contain no sensitive patient information. Moreover NRRD store the entire information in a single file as opposed to DICOM imaging.

2.2 Convert to a suitable format

As mentioned previously, NRRD provides a more insightful approach to understanding medical imaging and recognizing inherent patterns in a concised format. The conversion was done with the help of the Plastimatch tool which is an open source software for image computation. Plastimatch takes the DICOM image which is described in a polyline vectorized format, and converts it into a series of pixels which is more prominently known as rasterization. The

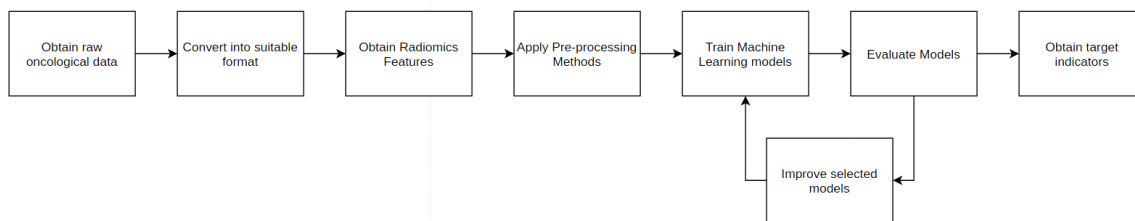


Fig. 1. Project Pipeline

TABLE 1
Data Description

Subtype	Number of patients	Estrogen Receptor	Progesterone Receptor	HER2	KI67 range
Luminal A	29	+	+/-	-	[5,20]
Luminal B	36	+	+/-	+/-	[25,80]
Triple Negative	19	-	-	-	[20,90]
HER	5	-	-	+	[30,50]

subroutine for rasterization of a DICOM image set with coordiantes x and y is shown below.

```
def rast(x, y, shape):
    nx, ny = draw.polygon(x, y, shape)
    nrrd = np.zeros(shape, dtype=np.bool)
    nrrd[ny, nx] = True

    return nrrd
```

Once this step is conducted, our image is in a compressed format, rife with information. Information extraction can be conducted through multiple means such as using neural networks, OCR recognition or pattern recognition algorithms.

2.3 Obtaining Radiomics Features

Information extraction from images directly has certain drawbacks. For eg, consider tumor classification using a standard Convolutional Neural Network (CNN). The CNN might be extremely successful in determining the existense of a blob of mass and it's exact location. However diagnosing the exact nature and feature set of the tumor is extremely difficult for a CNN. This is because a CNN views the image as simply a collection of pixels without any regard to the information embedded in all the views of the data.

To tackle this issue, we have utilized radiomics algorithms to extract feature sets from the medical images to reveal characteristics which are not captured by trained networks. The open-source Python library, PyRadiomics was used to mine out the required feature set. Before the actual extraction could be performed, a set of filters were applied on the NRRD to provide a comprehensive view of the data. The filters applied are listed in table 2.

To define a Region of Interest (ROI) and to check the dimensional constraints of the data, a mask file is utilized. The mask image corresponding to Figure 2 is shown in Figure 3. Note the red mark demarcating the tumor is done by a radiologist as is standard procedure. The features are now extracted from the image set with the help of the mask file. The features extracted are described by the Imaging

Biomarker Standardization Initiative (IBSI). The features have been shown in tables 3 and 4.

Therefore for each patient, the total number of features obtained are number of filters \times number of features i.e, $17 \times 100 = 1700$ features. Now that the entire feature set has been collected, we can begin the classification task.

2.4 Applying Pre-processing Techniques

From the 1700 features collected, not all of the features will contribute equally in the classification function. The process of preparing the input data for pattern learning by removing redundant characteristics, reducing noises and normalizing, selecting, and extracting features is termed as Data Pre-Processing.

3 CONCLUSION

The conclusion goes here.

APPENDIX A

PROOF OF THE FIRST ZONKLAR EQUATION

Appendix one text goes here.

APPENDIX B

Appendix two text goes here.

ACKNOWLEDGMENTS

The authors would like to thank...

REFERENCES

- [1] H. Kopka and P. W. Daly, *A Guide to L^AT_EX*, 3rd ed. Harlow, England: Addison-Wesley, 1999.



Fig. 2. Axial view with the tumor

TABLE 2
Applied Filters

Filter	Description	Equation
Wavelet (9)	Selective emphasizing or de-emphasizing of image in selected spatial frequency domain	-
Square	Square the image intensities	$x := (cx)^2$
Square Root	Compute root of image intensities	$x := \sqrt{cx}$
Laplacian of Gaussian $\sigma = 1, 2, 3$	Applies a Laplacian of Gaussian filter to the input image and yields a derived image for each sigma value specified	$\frac{1}{(\sigma\sqrt{2\pi})^3} e^{-\frac{x^2+y^2+z^2}{2\sigma^2}}$
Logarithm	Computes the natural logarithm of image intensities	$\text{clog}(x + 1)$
Exponential	Computes the exponential of the original image	e^{cx}
Gradient	Computes the gradient of the image	-

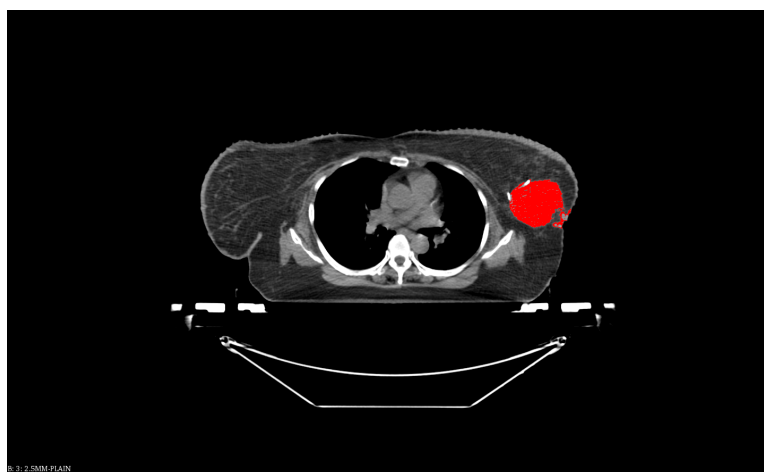


Fig. 3. Mask image

TABLE 3
Features-I

Feature Class	Feature	Feature Class	Feature
Shape	Elongation Flatness Least_Axis_length Major_Axis_Length Max_2D_Diameter_Column Max_2D_Diameter_Row Max_2D_Diameter_Slice Max_3D_Diameter Mesh_Volume Minor_Axis_Length Sphercity Surface_Area Surface_Volume Voxel_Volume	Grey Level Co-occurrence Matrix	Autocorrelation Cluster_Prominence Cluster_Shade Cluster_Tendency Constrast Correlation Difference_Average Difference_Entropy Difference_Variance Inverse_Variance Joint_Average Joint_Energy Joint_Entropy MCC Maximum_Probability Sum_Average Sum_Entropy Sum_Squares Id Idm Idn Idmn Imc1 Imc2
First Order Statistics	10 Percentile 90 Percentile Energy Entropy Interquartile_Range Kurtosis Maximum Mean_Absolute_Deviation Mean Median Minimum Range Robust_Mean_Deviation Robust_Mean_Squared Skewness Total_Energy Uniformity Variance	Grey Level Run Length Matrix	Uniformity Normalized_Uniformity Variance High_Run_Emphasis Long_Run_Emphasis Long_High_Run_Emphasis Long_Low_Run_Emphasis Low_Run_Emphasis Run_Entropy Run_Uniformity Run_Uniformity_Normalized Run_Percentage Run_Variance Short_Run_Emphasis Short_Run_High_Emphasis Short_Run_Low_Emphasis

TABLE 4
Features-II

Feature Class	Feature
Grey Level Size Zone Matrix	Non_Uniformity Non_Uniformity_Normalized Variance High_Zone_Emphasis Large_Area_Emphasis Large_Area_High_Level_Emphasis Large_Area_Low_Level_Emphasis Low_Zone_Emphasis Zone_Non_Uniformity Zone_Non_Uniformity_Normalized Small_Area_Emphasis Small_Area_High_Level_Emphasis Small_Area_Low_Level_Emphasis Zone_Entropy Zone_Percentage Zone_Variance
Gray Level Size Zone Matrix	Dependence_Entropy Dependence_Non_Uniformity Dependence_Non_Uniformity_Normalized Dependence_Variance GL_Non_Uniformity GL_Variance High_Emphasis Large_Dependence_Emphasis Large_Dependence_High_Emphasis Large_Dependence_Low_Emphasis Low_Emphasis Small_Dependence_Emphasis Small_Dependence_High_Emphasis Small_Dependence_Low_Emphasis
Neighbouring Gray Tone Difference Matrix	Busyness Coarseness Complexity Constrast Strength