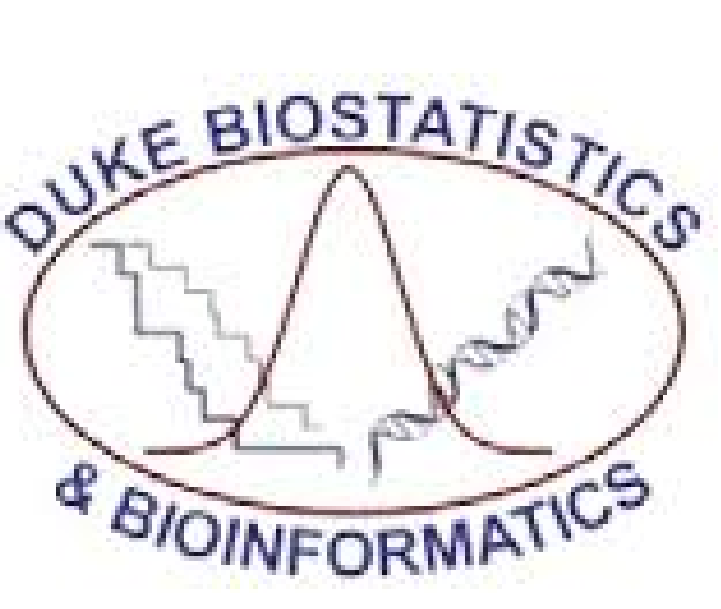




DISTINGUISH MALIGNANT FROM BENIGN BREAST CANCER

XIN LIU, JINGZHU ZHOU AND MENGSHU SHAO DEPARTMENT OF BIOSTATISTICS & BIOINFORMATICS



INTRODUCTION AND OBJECTIVES

Objective:The goal of this work is to use machine learning methodology to develop digitalized method to minimize subjectivity of the visual diagnosis and increase accuracy of classification for each breast cancer patient.
Introduction: Breast cancer stage can be diagnosed by visually examining fine needle aspirate

cancer cell images. The accuracy of the diagnosis is over 90%. However, there exists large standard error for the mean sensitivity and the mean specificity of the diagnosis. This indicates the accuracy varies greatly within individual series and visual diagnosis involves a great deal of subjectivity.

MATERIALS & METHODS

The following classification methods were considered in this research:

- Logistic regression
- Penalized logistic regression (LASSO plus Ridge)
- Linear / Quadratic Discriminant Analysis
- Support Vector Machine

The dataset contains records of 569 breast cancer instances, in which:

- 212 malignant
- 357 benign

Thirty variables are computed for each cell nucleus for breast cancer diagnosis prediction including mean(_m), sd(_sd), mean of largest three(_w) for the following features:

- 1) radius
- 2) texture
- 3) perimeter
- 4) area
- 5) smoothness
- 6) compactness
- 7) concavity
- 8) concave points
- 9) symmetry
- 10) fractal dimension



RESULTS 1

Model Selection for LDA,QDA,SVM:

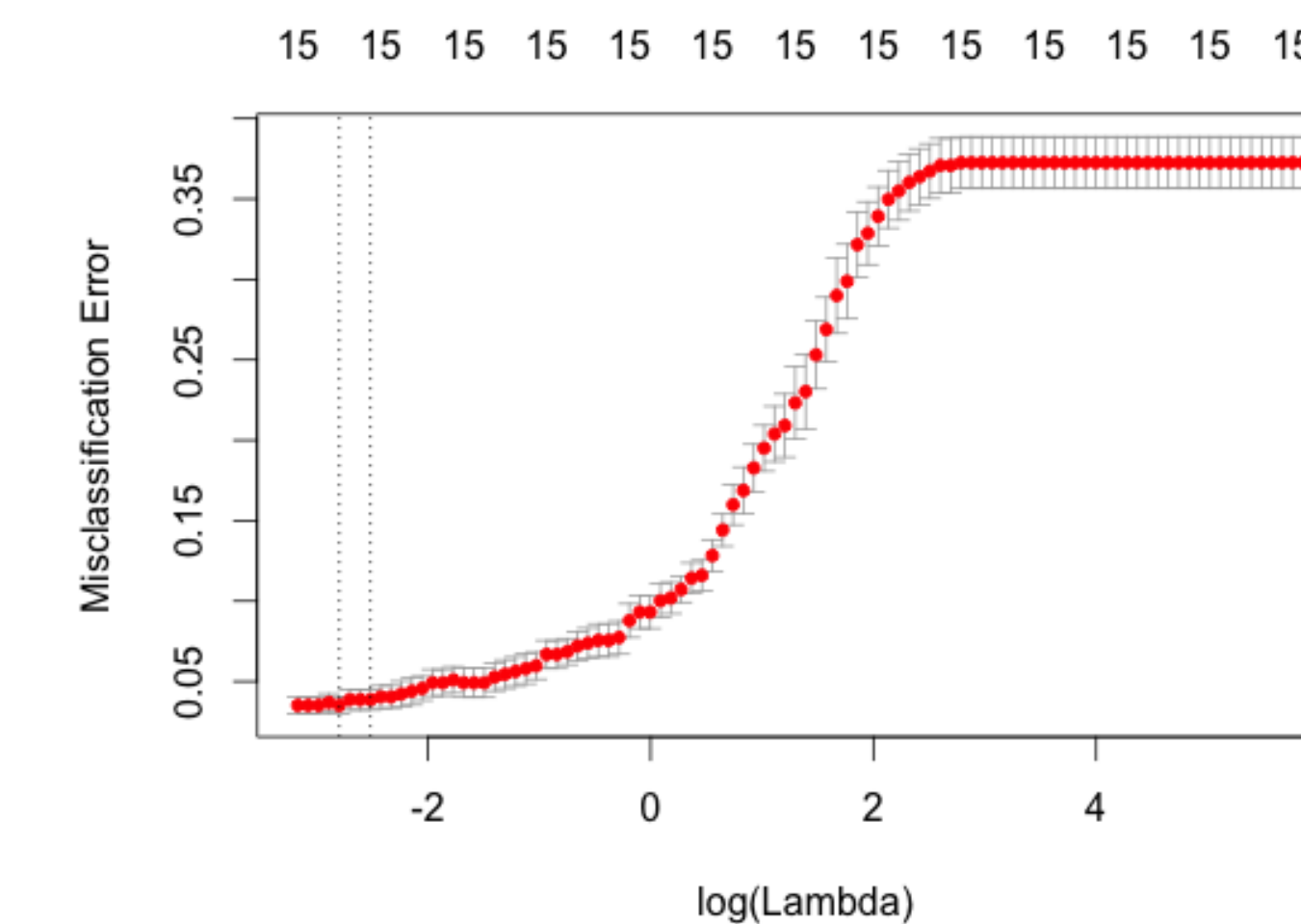
- 1.Pairwise correlations of the 30 variables were used to delete the variables that are highly correlated. Variable deleted: radius_m,perimeter_m,area_w,perimeter_w, radius_w,area_sd,perimeter_sd
- 2.Stepwise selection(logistic regression) were used to select the best subset of variables based on lowest AIC. Variable selected: concave_points_w, area_m, texture_w,radius_sd,compactness_sd, smoothness_w,concavity_w ,texture_sd

Model Selection for Penalized Regression:

- Lasso uses the constraint of the overall magnitude of the coefficients, thus important predictors are included in the model, and less important predictors shrink, potentially to zero. 10-folds cross validation selects the tuning parameter which makes the smallest cross validation error.
- The variables selected are texture_m, concavity_m, concave_points_m, fractal_m, radius_sd, texture_sd, smoothness_sd, compactness_sd, fractal_sd, radius_w, texture_w, smoothness_w, concavity_w, concave_points_w and symmetry_w.

RESULTS 2

1. Ridge parameter selection



The best tuning parameter:0.00304. Then Ridge regression was used to fit the model based on the tuning parameter.

$$\text{logit}(P(Y_i = 1|X = x)) = -9.35 + 0.073\text{texture}_m + 4.74\text{concavity}_m + 14.72\text{concave_points}_m - 39.34\text{fractal}_m + 1.81\text{radius}_sd - 0.13\text{texture}_sd - 8.03\text{smoothness}_sd - 2.52\text{compactness}_sd - 56.07\text{fractal}_sd + 0.15\text{radius}_w + 0.06\text{texture}_w + 12.44\text{smoothness}_w + 1.58\text{concavity}_w + 8.59\text{concave_points}_w + 4.41\text{symmetry}_w$$

2.Varying SVM input parameters

In another set of experiments we studied the dependence of the training error found by the SVM algorithm on the parameter C. Table 2 shows the variation of the training error corresponds with the parameter C. It can be seen that the smallest error is obtained for C=0.9.

Cost	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
Accuracy	93.85	95.08	96.31	96.49	96.84	97.19	97.19	97.54	97.72	97.36

Table 1: Parameter selection for SVM

CONCLUSION

Model	Logistic	Ridge	LDA	QDA	SVM
Accuracy	96.14	96.49	95.96	94.91	97.72
Predictor numbers	8	15	8	8	24

Table 2: Model Comparison

- All models give an accuracy rate above 95 %. Support vector machines gave the best accuracy, however it used the highest number of parameters.
- We would recommend logistic regression since it gave comparably high accuracy with the lowest number of parameters.
- Future Research:The various fitted models can be tested for accuracy, sensitivity and specificity on testing data (future patients prognosis) based on those cell measurements.Other classification methods such as extended bases LDA, QDA, plus interaction terms could be considered.

REFERENCES

- [1] William H. Wolberg, W. Nick Street, and Olvi L. Mangasarian. Computerized diagnosis of breast fine-needle aspirates. *The Breast Journal*, 3(2):77–80, 1997.
- [2] William H Wolberg, W.Nick Street, Dennis M Heisey, and Olvi L Mangasarian. Computer-derived nuclear features distinguish malignant from benign breast cytology. *Human Pathology*, 26(7):792 – 796, 1995.
- [3] William H Wolberg, W Nick Street, and Olvi L Mangasarian. Image analysis and machine learning applied to breast cancer diagnosis and prognosis. *Analytical and Quantitative cytology and histology*, 17(2):77–87, 1995.

CONTACT INFORMATION

Email xin.liu1@duke.edu
mengshu.shao@duke.edu
jingzhu.zhou@duke.edu