

Advanced Statistics: Application of supervised and unsupervised methods to biological data

Daniel Hill

2025-04-27

Abstract: Biological data with twenty features and four categorical class labels is explored and analysed using advanced statistical techniques in this report. Both supervised and unsupervised methods were implemented and evaluated, and the broad selection of models includes logistic regression, support vector machine, random forest, agglomerative hierarchical clustering, and gaussian mixture modelling. Comparing the results achieved using a selection of models with different underlying principles gives insight into the nature of the data. For example, the success of model-based clustering compared to hierarchical clustering and tree-based learning suggests the lack of hierarchy among the categorical labels, and the success of factor analysis as a dimensionality reduction technique suggests the presence of underlying biological mechanisms leading to several of the features arising together. Models achieving over 90% accuracy were produced, but all models performed notably worse at separating one of the categories that overlapped the other three.

Contents

1	Introduction	3
2	Methods	3
2.1	Data Description	3
2.2	Exploratory Data Analysis Approach	3
2.3	Supervised Learning Methods	6
2.4	Unsupervised Learning Methods	6
3	Results	6
3.1	Exploratory Data Analysis Findings	6
3.2	Supervised Learning Results	7
3.3	Unsupervised Learning Results	7
4	Discussion	7
5	Conclusion	7
6	References	7

List of Figures

1	Feature and Class Correlation Matrix: <i>highlighting relationships between variables and relationships with catagorical labels</i>	4
2	Distributions within feature columns: <i>histograms showing scale and skewness of data</i>	5
3	Distributions within feature columns: <i>Boxplots by class label by feature</i> . . .	5
4	PCA plotting of class labels: <i>scatterplot showing clustering tendency of catagorical classes</i>	6

List of Tables

1	Feature Descriptions	3
---	---------------------------------------	---

1 Introduction

this is a section where i write the introduction.

2 Methods

2.1 Data Description

20 features, a label with four catagorical classes. two groups of correlated features. Outliers removed using z score method. one feature transformed using logarithm. All features scaled and centered. All features numeric.

Bootstrapping was used to create a larger dataset.

2.2 Exploratory Data Analysis Approach

find distributions within each feature, look for correlations between features, scatter between plots that features that have high correlation to the catagorical label or another feature, use PCA to visualize all data together, calculate hopkins statistic to determine the clustering tendency of the data

Table 1: Feature Descriptions

Variable Name	No. missing values	mean	Std deviation	min	25th %ile	median	75th %ile	max
X1	2	9.876	0.764	6.840	9.356	9.872	10.399	12.355
X2	0	10.151	1.040	6.538	9.445	10.138	10.855	14.021
X3	1	8.861	0.871	6.424	8.243	8.847	9.445	12.216
X4	2	8.939	1.275	3.875	8.088	8.926	9.805	13.351
X5	0	13.853	0.942	10.527	13.236	13.858	14.492	16.557
X6	0	8.151	1.026	4.815	7.447	8.134	8.856	11.871
X7	1	0.426	0.278	0.000	0.185	0.375	0.678	1.301
X8	3	0.234	0.197	0.000	0.102	0.170	0.300	1.230
X9	0	0.717	0.247	0.006	0.532	0.751	0.888	1.679
X10	0	0.378	0.155	0.000	0.281	0.372	0.469	1.199
X11	1	9.175	1.087	6.031	8.412	9.077	9.860	13.027
X12	0	11.930	0.977	8.046	11.290	11.913	12.594	15.478
X13	1	8.228	0.806	4.919	7.719	8.244	8.746	11.226
X14	2	7.846	1.238	3.574	7.022	7.884	8.689	12.413
X15	1	10.701	0.962	7.572	10.054	10.701	11.344	14.037
X16	1	7.814	1.052	3.801	7.132	7.830	8.521	11.668
X17	2	0.504	0.221	0.001	0.348	0.502	0.653	1.315
X18	3	0.682	0.204	0.004	0.544	0.686	0.819	1.390
X19	0	0.544	0.254	0.000	0.363	0.545	0.710	1.518
X20	0	0.589	0.231	0.012	0.434	0.587	0.746	1.353

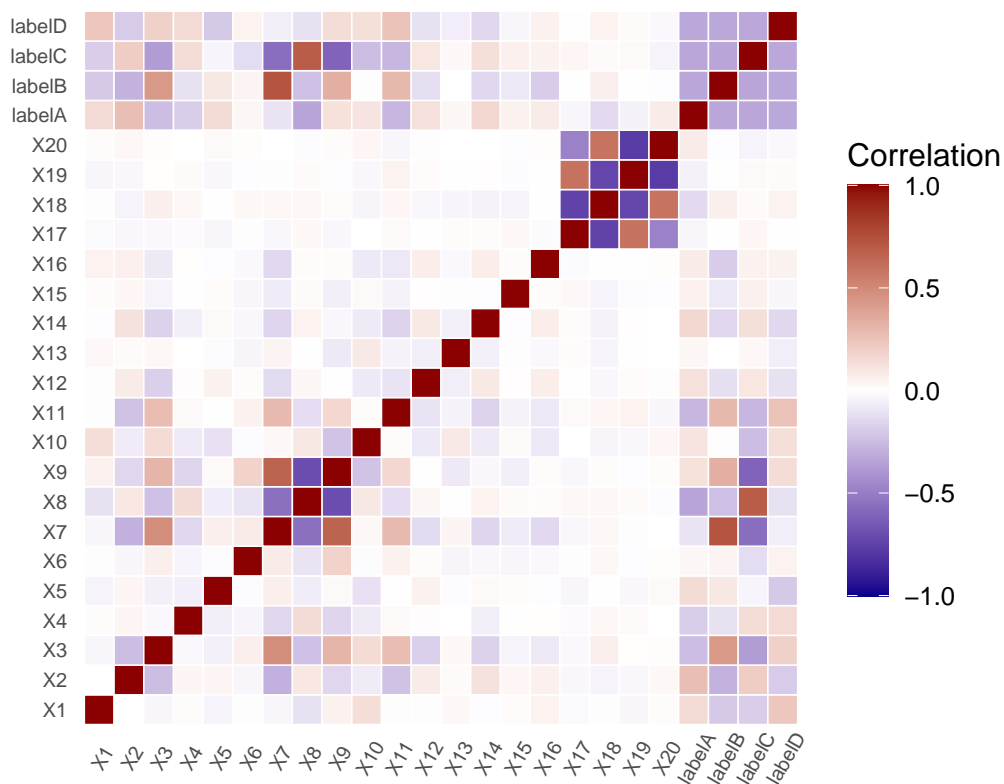


Figure 1: **Feature and Class Correlation Matrix:** *highlighting relationships between variables and relationships with catagorical labels*

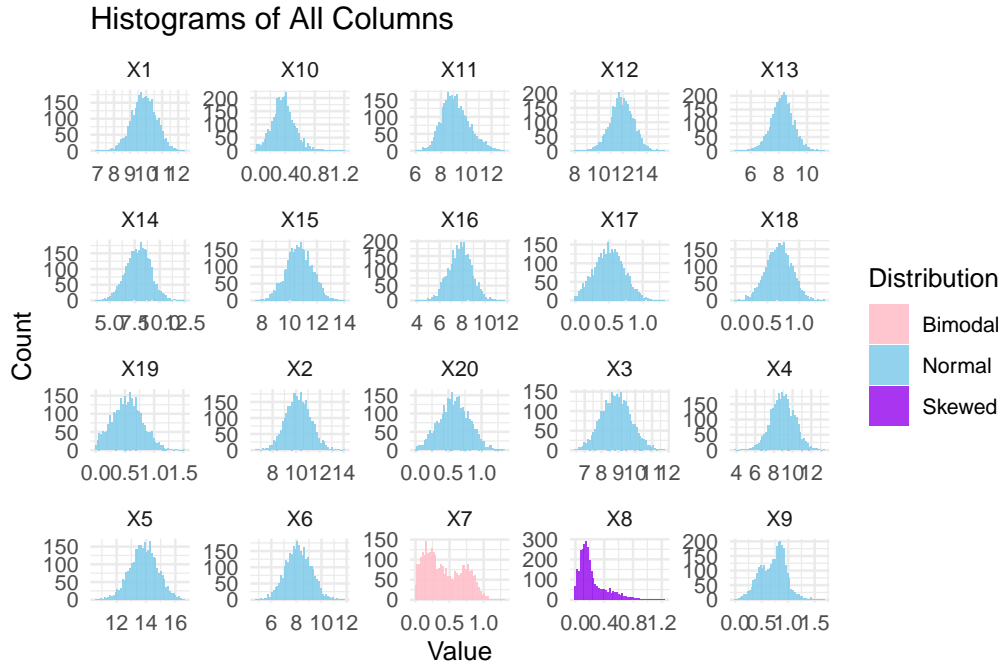


Figure 2: **Distributions within feature columns:** *histograms showing scale and skewness of data*

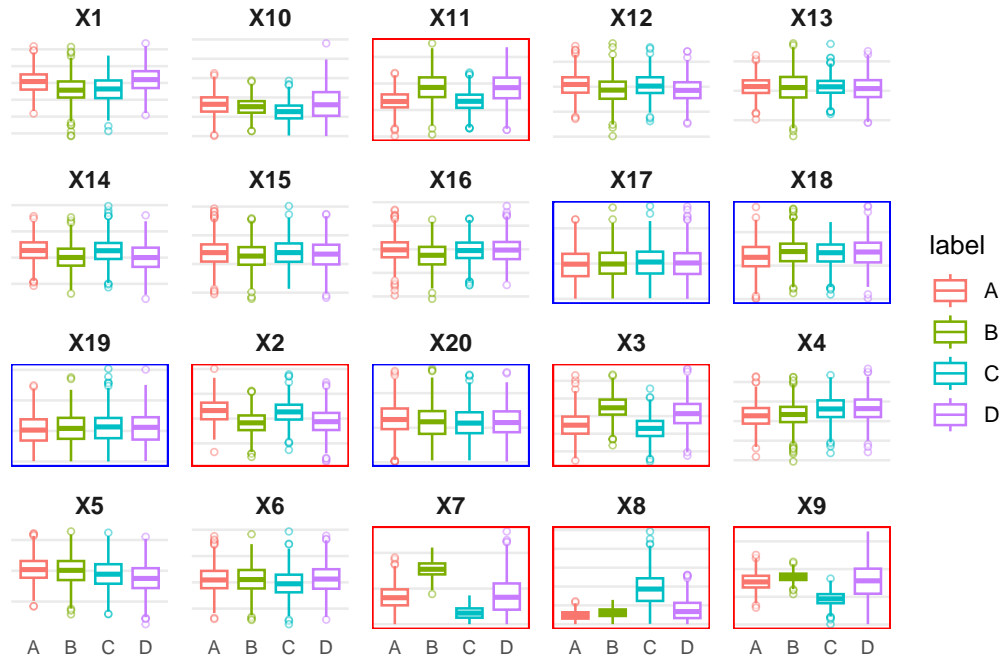


Figure 3: **Distributions within feature columns:** *Boxplots by class label by feature*

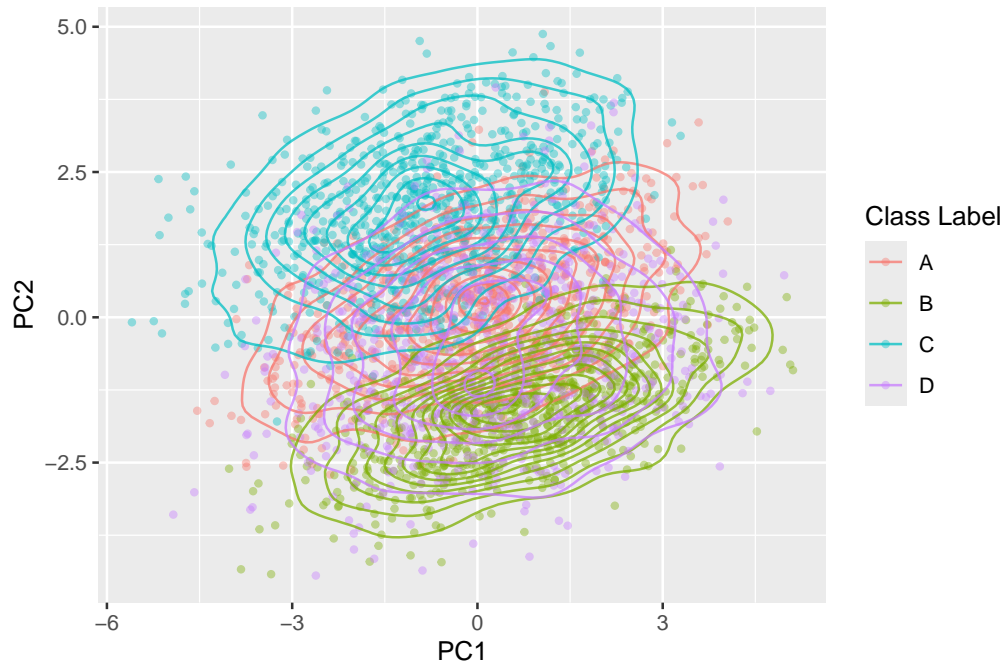


Figure 4: **PCA plotting of class labels:** *scatterplot showing clustering tendency of catagorical classes*

2.3 Supervised Learning Methods

logistic regression - including class weighting and L2 regularization, and feature selection. Random forest, including feature selection and tuning of mtry. SVM, with feature selection and tuning of kernel selection, gamma and cost parameters.

2.4 Unsupervised Learning Methods

Agglomerative hierarchical clustering, including tuning of linking metric. Gaussian mixture model based clustering, including selecting a model, regularization using shrinkage parameter, and dimensionality reduction using factor analysis.

3 Results

3.1 Exploratory Data Analysis Findings

most features are normally distributed except for X8 which is highly skewed. Features originally had different scales. X7, X8, X9 columns are correlated and correlate highly with the labels. X17, X18, X19 and X20 are highly correlated together and have very low correlation with the labels.

The PCA showed the four labels had some clustered structure, but also some significant overlap. The hopkins statistic showed there was a moderate clustering tendency, but that the label column when included made the clustering tendency extremely high. This is an initial suggestion that supervised learning would be more effective than unsupervised learning

3.2 Supervised Learning Results

logistic regression was not great. weighting did nothing, as expected. regularization didn't really do anything. feature selection did not improve the model. We saw that the model particularly underperformed at classifying class D correctly.

The random forest performed well without any configuration. feature selection was not effective. still struggled at separating class D. mtry was tuned.

svm was good but not as good as random forest. lots of tuning. feature selection was ineffective

3.3 Unsupervised Learning Results

ahc was good not great.

Gaussian mixed model clustering performed very well. dimensionality reduction using factor analysis was slightly effective.

4 Discussion

the final model had good overall accuracy but caution is advised when using a ml model with this data due to the poor performance in class D - if false positives or false negatives in this class have serious implications, some models become immediately unusable.

5 Conclusion

6 References