

# Report of data analysis for early-stage epithelial ovarian cancer and the comparison to the previous results thereof

Fan Zhang

## 1. Background

This analysis has heretofore intended to repeat the previous study from Joselle O'Brien regarding the predictive modeling for early-stage epithelial ovarian cancer using CA125, Mass Spectrometry (MS) and Nuclear Magnetic Resonance (NMR) Spectroscopy data. The analysis started from pre-processing of the raw data, followed by modeling using kNN, Partial Least Square (PLS), Bagging, RandomForest (RF), Support Vector Machine (SVM), Logistic regression and Boosting for each dataset and each classification model, i.e., Healthy and Benign lumped (HB) model, Malignant and Benign lumped (MB) model, original three classes model, respectively.

## 2. Methods and Results

### 2.1 Pre-processing

Variables from MS dataset have been filtered such that those variables of which values are only 1's are deleted. By checking with Joselle's dataset, observations of g214 and g275 were also deleted from the MS dataset. Each dataset was separated into two parts, XX and YY, with XX representing all dependent variables and YY being one dimensional vector indicating the diagnostic classification. For MS data, each variable's name has been replaced by the associated compound ID prefixed by an "X". Details of each dataset are listed below in Table 1:

Dataset	Raw data		Clean data			
	XX	YY coding	XX	YY coding	HB coding	MB coding
CA-125	120 x 1	Benign Cancer Normal	120 x 1	"Normal"-1, "Benign"-2, "Cancer"-3	Normal/Benign-Negative-1 Cancer-Positive-2	Normal-Negative-1 Benign/Cancer-Positive-2
MS	120 x 758	Benign Cancer Healthy	118 x 720	"Healthy"-1, "Benign"-2, "Cancer"-3	Healthy/Benign-Negative-1 Cancer-Positive-2	Healthy-Negative-1 Benign/Cancer-Positive-2
DIRE	--	--	118 x 3055	"Normal"-1, "Benign"-2, "Cancer"-3	Normal/Benign-Negative-1 Cancer-Positive-2	Normal-Negative-1 Benign/Cancer-Positive-2
DOSY	--	--	118 x 3053	"Normal"-1, "Benign"-2, "Cancer"-3	Normal/Benign-Negative-1 Cancer-Positive-2	Normal-Negative-1 Benign/Cancer-Positive-2
CPMGc	--	--	118 x 3055	"Normal"-1, "Benign"-2, "Cancer"-3	Normal/Benign-Negative-1 Cancer-Positive-2	Normal-Negative-1 Benign/Cancer-Positive-2
CPMGm	--	--	118 x 3057	"Normal"-1, "Benign"-2, "Cancer"-3	Normal/Benign-Negative-1 Cancer-Positive-2	Normal-Negative-1 Benign/Cancer-Positive-2
NOESYc	--	--	118 x 3055	"Normal"-1,	Normal/Benign-Negative-1	Normal-Negative-1

			"Benign"-2, "Cancer"-3	Cancer-Positive-2	Benign/Cancer-Positive-2
NOESYm	--	--	118 x 3057 "Normal"-1, "Benign"-2, "Cancer"-3	Normal/Benign-Negative-1 Cancer-Positive-2	Normal-Negative-1 Benign/Cancer-Positive-2

**Table 1: Dimensions for each dataset**

## 2.2 Parameter scanning and model building

Principle component analysis (PCA) was not included in the modeling part, but an individual level plot (Figure 1) was created to give a glance over the classification of the observations for the MS data.



**Figure 1: Individuals factor map.** The X axis represents the first principle component (PC) and the Y axis represents the second PC, percentages in the parenthesis indicate the X variance explained by the corresponding PC, respectively. Each observation is projected onto the map according to the scores from both PCs. It can be seen that the Benign group is in the middle between Cancer and Healthy groups.

After PCA, modeling using kNN and logistic regression were performed on CA125 dataset. MS and NMR datasets were modeled by kNN, PLS, Bagging, RF, Boosting and SVM. When it's modeled by SVM, three different kernels, i.e., linear, polynomial and radial kernels, were tried with different parameters. Since all methods have different tuning parameters, wide ranges of parameters or combinations of parameters for different methods were searched for best predicative power. Due to computational limit, the ranges of some parameters searched for all six

NMR datasets were narrower than that of MS dataset (Table 2).

Method	Parameter	MS	NMR (6 sets)
kNN	K	1, 2, ..., 30	1, 2, ..., 15
PLS	Latent Variables	1, 2, ..., 30	1, 2, ..., 15
Bagging	# trees	3000	3000
RF	mtry	Sqrt(p)	Sqrt(p)
	# trees	3000	3000
Boosting	Depth	1, 2, ..., 5	1, 2, 3
	Shrinkage	0.01, 0.001	0.01, 0.001
	#trees	1000, 5000	10, 100
SVM	Linear	Cost	$2^{(-10)}, 2^{(-9)}, \dots, 2^{15}$
	Polynomial	Cost	$2^{(-10)}, 2^{(-9)}, \dots, 2^{15}$
		Degree	2, 3, 4, 5
	Radial	Cost	$2^{(-10)}, 2^{(-9)}, \dots, 2^{15}$
		gamma	$2^{(-15)}, 2^{(-14)}, \dots, 2^{15}$

Table 2: Parameters setting

Cross validation using three folds were performed to get the mean prediction accuracy. The highest mean accuracy of different methods among various parameter settings is listed in Table 3 with respect to all the datasets and in Table 4 with respect to classification models.

Method	Dataset								Model
	CA125	MS	CPMGc	CPMGm	DIRE	DOSY	NOESYc	NOESYm	
kNN	0.833/0.808	0.796/0.797	0.807/0.814	0.712/0.678	0.744/0.661	0.762/0.746	0.889/0.822	0.787/0.746	HB
	0.708/0.625	0.813/0.847	0.874/0.873	0.771/0.703	0.838/0.703	0.779/0.745	0.882/0.831	0.832/0.772	MB
	0.583/0.492	0.610/0.703	0.747/0.772	0.551/0.407	0.634/0.415	0.753/0.703	0.864/0.796	0.576/0.509	3-Class
PLS	--	0.737/0.822	0.789/0.813	0.644/0.754	0.721/0.754	0.88/0.865	0.846/0.898	0.78/0.746	HB
	--	0.804/0.856	0.873/0.872	0.805/0.823	0.788/0.797	0.897/0.889	0.871/0.898	0.806/0.822	MB
	--	0.484/0.652	0.61/0.576	0.542/0.483	0.552/0.475	0.695/0.593	0.625/0.643	0.517/0.593	3-Class
Bagging	--	0.814/0.839	0.747/0.746	0.652/0.695	0.693/0.679	0.814/0.822	0.856/0.847	0.72/0.754	HB
	--	0.837/0.873	0.804/0.856	0.73/0.771	0.762/0.746	0.779/0.804	0.839/0.873	0.788/0.780	MB
	--	0.754/0.763	0.746/0.712	0.525/0.551	0.599/0.653	0.702/0.720	0.814/0.881	0.603/0.602	3-Class
RF	--	0.814/0.805	0.721/0.762	0.669/0.703	0.668/0.713	0.814/0.797	0.873/0.847	0.754/0.771	HB
	--	0.864/0.831	0.804/0.856	0.737/0.771	0.721/0.754	0.753/0.788	0.858/0.839	0.771/0.780	MB
	--	0.771/0.746	0.763/0.746	0.485/0.525	0.6/0.653	0.736/0.788	0.882/0.872	0.577/0.593	3-Class
Boosting	--	0.881/0.848	0.823/0.703	0.72/0.687	0.719/0.687	0.831/0.720	0.924/0.839	0.771/0.721	HB
	--	0.855/0.890	0.88/0.762	0.848/0.686	0.788/0.669	0.821/0.695	0.89/0.856	0.865/0.797	MB
	--	0.814/0.840	0.805/0.559	0.595/0.407	0.65/0.534	0.737/0.619	0.847/0.644	0.618/0.399	3-Class
SVM-L	--	0.763/0.797	0.814/0.78	0.627/0.703	0.779/0.789	0.839/0.856	0.872/0.839	0.745/0.737	HB
	--	0.796/0.805	0.84/0.865	0.771/0.797	0.846/0.865	0.873/0.813	0.898/0.864	0.815/0.772	MB
	--	--	--	--	--	--	--	--	3-Class
SVM-P	--	0.763/0.78	0.814/0.661	0.627/0.661	0.779/0.661	0.839/0.661	0.872/0.66	0.745/0.661	HB
	--	0.796/0.797	0.84/0.660	0.771/0.661	0.846/0.66	0.873/0.661	0.898/0.661	0.815/0.662	MB
	--	--	--	--	--	--	--	--	3-Class
SVM-R	--	0.763/0.78	0.814/0.822	0.627/0.703	0.779/0.789	0.839/0.890	0.872/0.898	0.745/0.788	HB
	--	0.796/0.839	0.84/0.865	0.771/0.797	0.846/0.898	0.873/0.847	0.898/0.915	0.815/0.78	MB
	--	--	--	--	--	--	--	--	3-Class
Logistic	0.817/0.833	--	--	--	--	--	--	--	HB
	0.642/0.617	--	--	--	--	--	--	--	MB

0.508/0.508	--	--	--	--	--	--	--	3-Class
-------------	----	----	----	----	----	----	----	---------

**Table 3: Mean accuracy for datasets. Values on the left side of "/" represent Joselle's results, values on the right side represent my result. Red color indicates increasing, green for decreasing, black for tie. Dashed box indicates the highest value in that column.**

Dataset	HB	MB	3-Class	Method	HB	MB	3-Class	Method
CA125	0.833/0.808	0.708/0.625	0.583/0.492		--	--	--	
MS	0.796/0.797	0.813/0.847	0.610/0.703		0.763/0.797	0.796/0.805	--	
CPMGc	0.807/0.814	0.874/0.873	0.747/0.772		0.814/0.78	0.84/0.865	--	
CPMGm	0.712/0.678	0.771/0.703	0.551/0.407	kNN	0.627/0.703	0.771/0.797	--	SVM-L
DIRE	0.744/0.661	0.838/0.703	0.634/0.415		0.779/0.789	0.846/0.865	--	
DOSY	0.762/0.746	0.779/0.745	0.753/0.703		0.839/0.856	0.873/0.813	--	
NOESYc	0.889/0.822	0.882/0.831	0.864/0.796		0.872/0.839	0.898/0.864	--	
NOESYm	0.787/0.746	0.832/0.772	0.576/0.509		0.745/0.737	0.815/0.772	--	
CA125	--	--	--		--	--	--	
MS	0.737/0.822	0.804/0.856	0.484/0.652		0.763/0.78	0.796/0.797	--	
CPMGc	0.789/0.813	0.873/0.872	0.61/0.576		0.814/0.661	0.84/0.660	--	
CPMGm	0.644/0.754	0.805/0.823	0.542/0.483	PLS	0.627/0.661	0.771/0.661	--	SVM-P
DIRE	0.721/0.754	0.788/0.797	0.552/0.475		0.779/0.661	0.846/0.66	--	
DOSY	0.88/0.865	0.897/0.889	0.695/0.593		0.839/0.661	0.873/0.661	--	
NOESYc	0.846/0.898	0.871/0.898	0.625/0.643		0.872/0.66	0.898/0.661	--	
NOESYm	0.78/0.746	0.806/0.822	0.517/0.593		0.745/0.661	0.815/0.662	--	
CA125	--	--	--		--	--	--	
MS	0.814/0.839	0.837/0.873	0.754/0.763		0.763/0.78	0.796/0.839	--	
CPMGc	0.747/0.746	0.804/0.856	0.746/0.712		0.814/0.822	0.84/0.865	--	
CPMGm	0.652/0.695	0.73/0.771	0.525/0.551	Bagging	0.627/0.703	0.771/0.797	--	SVM-R
DIRE	0.693/0.679	0.762/0.746	0.599/0.653		0.779/0.789	0.846/0.898	--	
DOSY	0.814/0.822	0.779/0.804	0.702/0.720		0.839/0.890	0.873/0.847	--	
NOESYc	0.856/0.847	0.839/0.873	0.814/0.881		0.872/0.898	0.898/0.915	--	
NOESYm	0.72/0.754	0.788/0.780	0.603/0.602		0.745/0.788	0.815/0.78	--	
CA125	--	--	--		--	--	--	
MS	0.814/0.805	0.864/0.831	0.771/0.746		0.881/0.848	0.855/0.890	0.814/0.840	
CPMGc	0.721/0.762	0.804/0.856	0.763/0.746		0.823/0.703	0.88/0.762	0.805/0.559	
CPMGm	0.669/0.703	0.737/0.771	0.485/0.525	RF	0.72/0.687	0.848/0.686	0.595/0.407	Boosting
DIRE	0.668/0.713	0.721/0.754	0.6/0.653		0.719/0.687	0.788/0.669	0.65/0.534	
DOSY	0.814/0.797	0.753/0.788	0.736/0.788		0.831/0.720	0.821/0.695	0.737/0.619	
NOESYc	0.873/0.847	0.858/0.839	0.882/0.872		0.924/0.839	0.89/0.856	0.847/0.644	
NOESYm	0.754/0.771	0.771/0.780	0.577/0.593		0.771/0.721	0.865/0.797	0.618/0.399	
CA125	0.817/0.833	0.642/0.617	0.508/0.508					
MS	--	--	--					
CPMGc	--	--	--					
CPMGm	--	--	--	Logistic				
DIRE	--	--	--					
DOSY	--	--	--					
NOESYc	--	--	--					
NOESYm	--	--	--					

**Table 4: Mean accuracy for classification models. Values on the left side of "/" represent Joselle's results, values on the right side represent my result. Red color indicates increasing, green for decreasing, black for tie. Dashed**

**box indicates the highest value in that column.**

Table 3 and 4 are of same values but in different arrangements. From table 3, we can see that the highest accuracy for CA125 data is 83.3% from both Joselle's and my result, and each method from kNN and logistic regression can achieve it. For MS data, it's 89% as opposed to 88.1%, and both were obtained by boosting. For NMR data, the highest accuracy occurs in NOESYc dataset achieved by SVM-Radial, it's 91.5% while the previous one is 92.4% obtained by boosting. In regards to classification models (Table 4), the highest accuracies from HB and MB models are 89.8% and 91.5%, both obtained by SVM-Radial. But the accuracy for 3-class model is slightly lower, which is 88.1% as opposed to 88.2%, from Bagging and RF, respectively.

If takes standard deviation into account as the accuracy was obtained from three folds (Figure 2-4), it would be clear that SVM stands out as the best choice for NMR data, i.e., NOESYc, in terms of HB and MB models while the Bagging is the best for 3-Class model thereof. For MS data, it turned out that the Boosting is uniformly superior among all methods for all three models. Taking account of highest mean accuracy and smallest SD, the models chosen and the associated parameters are concluded in Table 5.

Data source	Chosen method	Parameters	Model
CA125	kNN	K=17	HB
	kNN	K=25	MB
	kNN	K=29	3-class
MS	Boosting	TreeNum=1000, Shrinkage=0.01, Depth=2	HB
	Boosting	TreeNum=5000, Shrinkage=0.01, Depth=2	MB
	Boosting	TreeNum=5000, Shrinkage=0.001, Depth=1	3-class
NMR (NOESYc)	SVM Radial	Cost=32, Gamma=1024	HB
	SVM Radial	Cost=128, Gamma=256	MB
	Bagging	TreeNum=3000	3-class

**Table 5: Chosen models and the corresponding parameters.**

## 2.3 Consensus results

Consensus analysis was performed based on the chosen methods from last section for each model. Those selected models were refit to each dataset, i.e., CA125, MS and NOESYc, and their predictions for each subject were merged together, after which a majority vote was conducted in a way that the final prediction for each subject was the majority vote among the three datasets. The results are in Table 6-8. In conclusion, the accuracies are 85%, 90% and 82% for HB, MB and 3-class models, respectively.

## Appendix

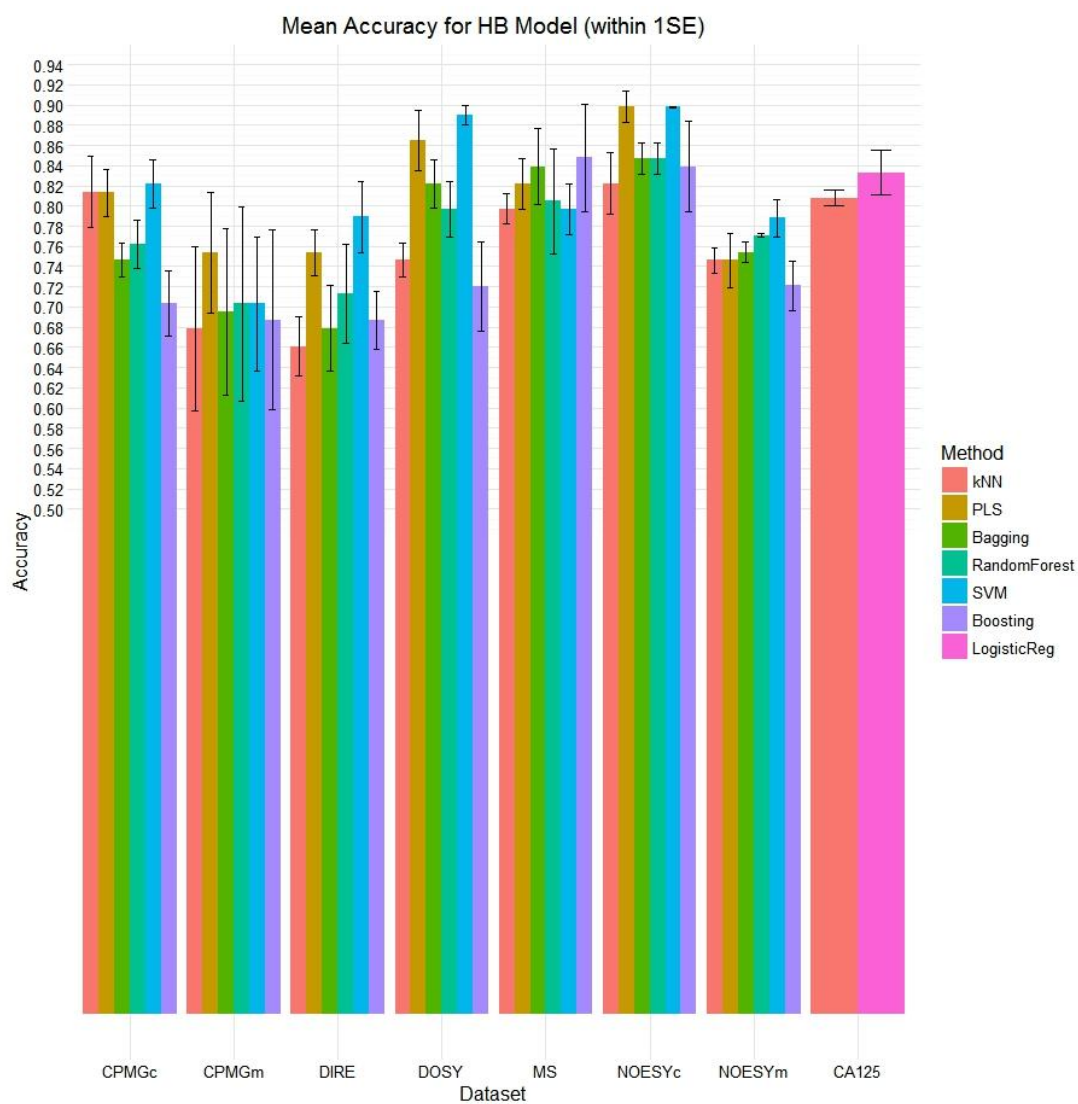


Figure 2: Mean accuracy comparison for HB model.

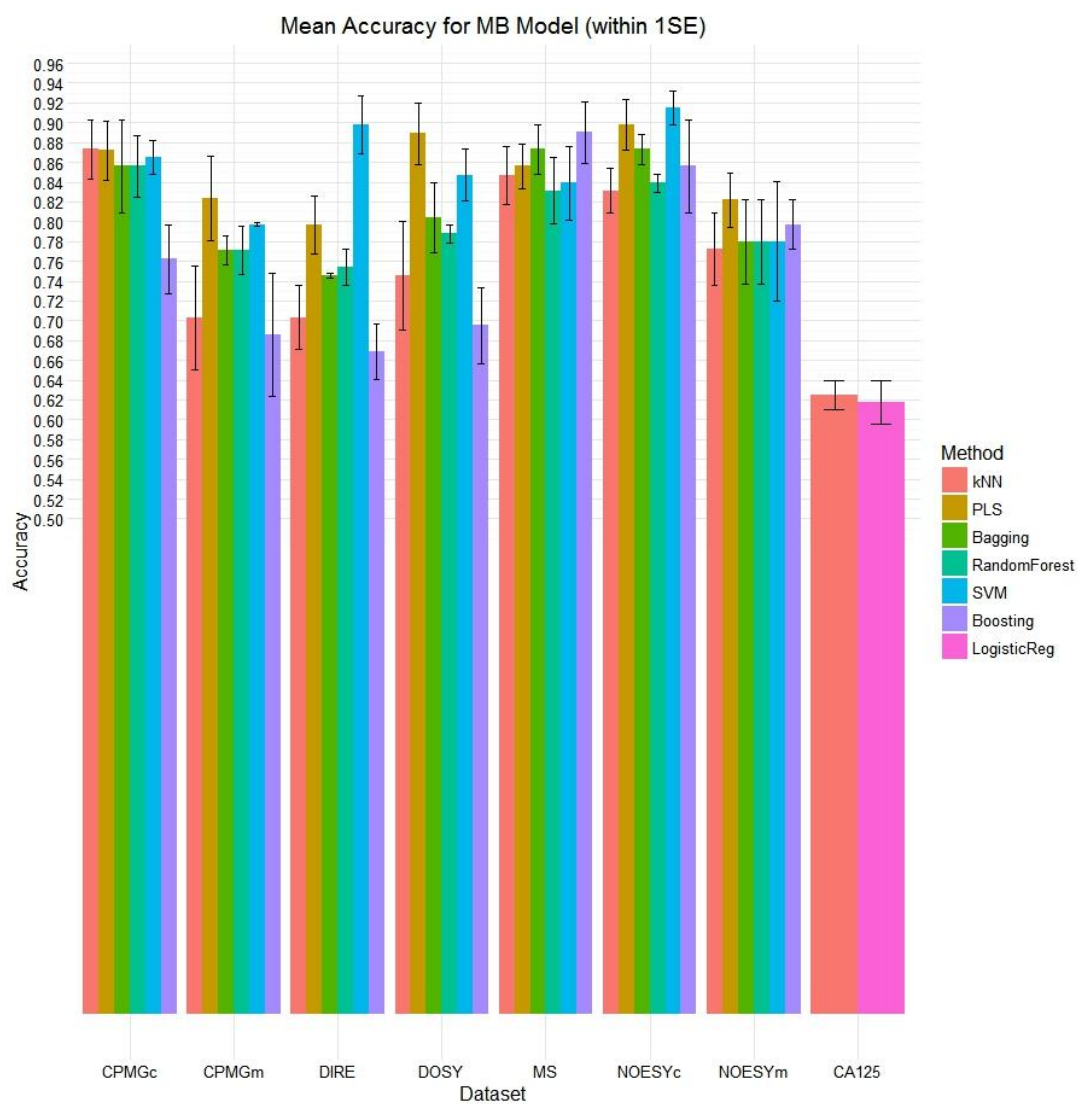


Figure 3: Mean accuracy comparison for MB model.

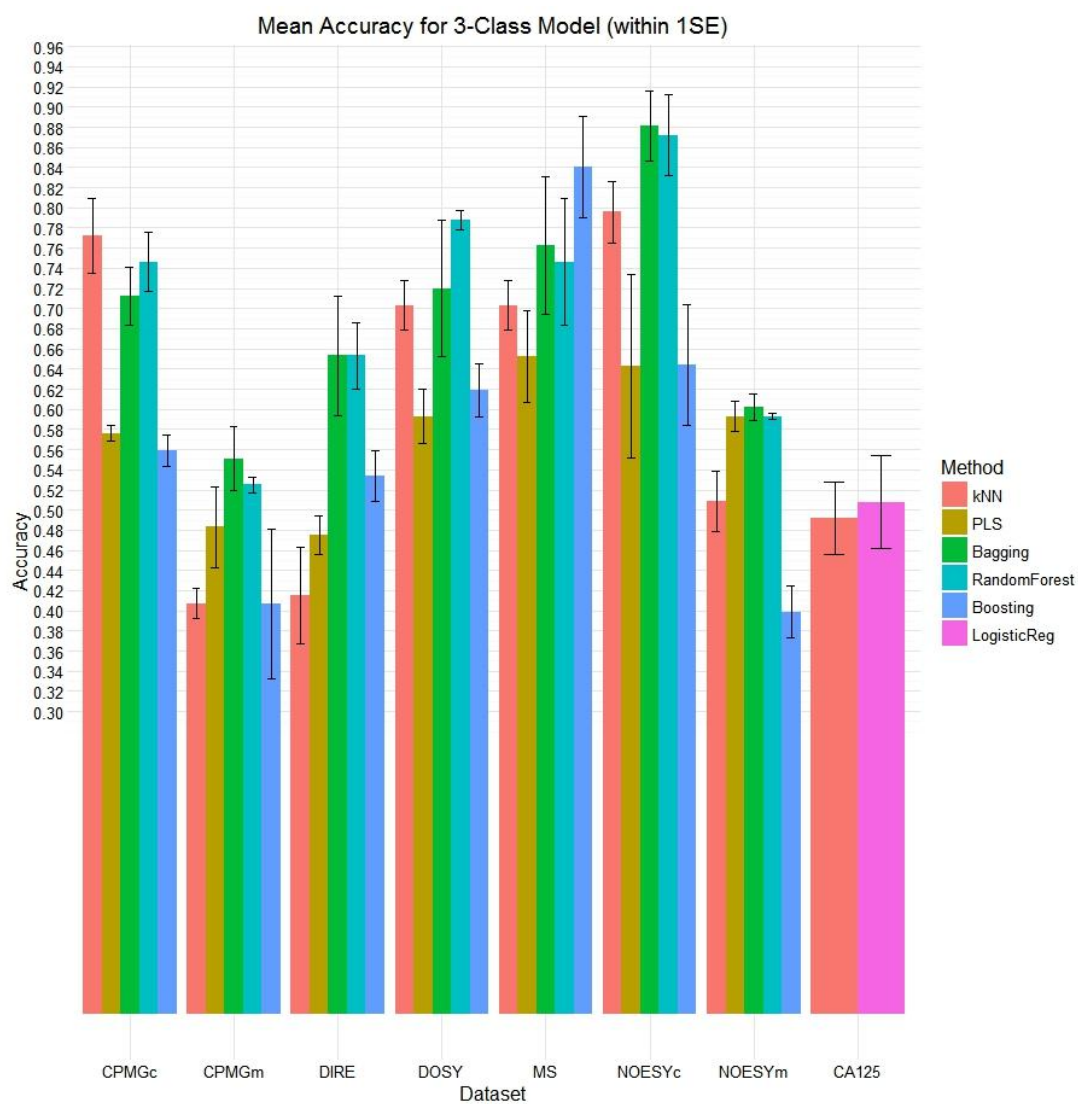


Figure 4: Mean accuracy comparison for 3-Class model.



ID	CA125	MS	NMR	Vote	TRUE	ID	CA125	MS	NMR	Vote	TRUE	ID	CA125	MS	NMR	Vote	TRUE
g111	1	1	1	1	1	g208	1	1	1	1	1	g275	2	NA	1	99	1
g113	1	1	1	1	1	g209	1	1	1	1	1	g276	1	1	1	1	1
g119	1	1	1	1	1	g21	1	1	1	1	2	g277	1	1	1	1	1
g122	1	1	1	1	1	g210	1	1	2	1	1	g28	1	2	2	2	2
g124	1	1	2	1	1	g213	1	1	1	1	1	g281	1	1	1	1	1
g126	1	1	1	1	1	g214	1	NA	1	1	1	g286	1	1	1	1	1
g129	1	1	1	1	1	g215	1	1	1	1	1	g289	1	1	1	1	1
g135	2	2	2	2	2	g216	1	1	1	1	1	g291	1	1	1	1	1
g139	1	1	1	1	2	g217	1	1	1	1	1	g293	1	1	1	1	1
g140	2	1	1	1	2	g220	1	1	1	1	1	g295	1	1	1	1	1
g142	1	1	2	1	2	g221	1	1	1	1	1	g296	1	1	1	1	1
g143	1	1	1	1	2	g223	1	1	1	1	1	g298	1	1	1	1	1
g144	1	1	2	1	2	g225	1	1	1	1	1	g301	1	1	1	1	1
g147	2	2	2	2	2	g227	1	1	1	1	1	g303	1	1	1	1	1
g149	1	2	2	2	2	g23	1	1	2	1	2	g305	1	1	1	1	1
g15	1	2	1	1	2	g230	1	1	1	1	1	g306	1	1	1	1	1
g150	2	2	2	2	2	g233	1	1	1	1	1	g307	1	1	1	1	1
g151	2	2	2	2	2	g234	1	1	1	1	1	g308	1	1	1	1	1
g152	1	1	2	1	2	g237	1	1	1	1	1	g309	1	1	1	1	1
g155	2	2	2	2	2	g239	1	1	1	1	1	g314	1	1	1	1	1
g156	1	2	2	2	2	g241	1	1	1	1	1	g318	1	1	1	1	1
g159	1	1	1	1	2	g242	1	1	1	1	1	g326	1	1	1	1	1
g160	1	2	2	2	2	g244	1	1	1	1	1	g333	1	2	NA	99	1
g162	2	1	2	2	2	g245	1	1	1	1	1	g337	1	1	1	1	1
g17	1	2	1	1	2	g249	1	1	1	1	1	g339	1	1	1	1	1
g170	2	1	1	1	2	g250	1	1	1	1	1	g340	1	1	1	1	1
g173	2	2	2	2	2	g252	1	1	2	1	2	g342	1	1	1	1	1
g176	2	2	2	2	2	g256	1	1	NA	1	1	g343	1	1	1	1	1
g178	2	2	2	2	2	g257	1	1	1	1	1	g38	2	2	1	2	2
g179	1	1	2	1	2	g259	2	1	1	1	1	g40	2	2	2	2	2
g181	2	2	1	2	2	g26	2	2	2	2	2	g44	2	2	1	2	2
g188	2	2	1	2	2	g261	1	1	1	1	1	g48	2	1	1	1	2
g19	1	2	2	2	2	g265	1	1	1	1	1	g7	1	1	1	1	2
g191	2	2	2	2	2	g266	1	1	1	1	1	g72	1	1	1	1	1
g192	2	2	2	2	2	g267	1	2	1	1	1	g73	1	1	2	1	1
g193	1	2	2	2	2	g269	1	1	1	1	1	g74	1	1	1	1	1
g199	1	1	1	1	1	g270	1	2	1	1	1	g75	1	1	1	1	1
g201	1	1	1	1	1	g271	1	1	1	1	1	g77	1	2	1	1	1
g203	1	1	1	1	1	g272	1	1	1	1	1	g8	1	2	2	2	2
g205	1	1	1	1	1	g274	2	1	1	1	1	g80	1	1	1	1	1

Table 6: HB consensus results. A majority vote value of "99" indicates that the predictions from the three datasets are all different on this subject.

ID	CA125	MS	NMR	Vote	TRUE	ID	CA125	MS	NMR	Vote	TRUE	ID	CA125	MS	NMR	Vote	TRUE
g111	2	1	2	2	1	g208	2	1	1	1	1	g275	2	NA	2	2	2
g113	2	2	1	2	1	g209	2	1	1	1	1	g276	2	2	2	2	2
g119	1	1	1	1	1	g21	2	2	1	2	2	g277	1	2	2	2	2
g122	1	1	1	1	1	g210	1	2	2	2	1	g28	2	2	2	2	2
g124	1	2	1	1	1	g213	2	1	1	1	1	g281	1	2	2	2	2
g126	1	2	1	1	1	g214	1	NA	1	1	1	g286	2	2	2	2	2
g129	2	2	1	2	1	g215	2	1	1	1	1	g289	2	2	2	2	2
g135	2	2	2	2	2	g216	1	1	1	1	1	g291	2	2	2	2	2
g139	2	2	1	2	2	g217	1	1	1	1	1	g293	1	2	2	2	2
g140	2	1	1	1	2	g220	1	1	1	1	1	g295	1	2	2	2	2
g142	2	2	2	2	2	g221	2	1	1	1	1	g296	1	2	2	2	2
g143	2	1	1	1	2	g223	2	1	1	1	1	g298	1	2	2	2	2
g144	2	2	2	2	2	g225	1	1	1	1	1	g301	2	2	2	2	2
g147	2	2	2	2	2	g227	1	1	1	1	1	g303	2	2	2	2	2
g149	2	2	2	2	2	g23	1	2	2	2	2	g305	2	2	2	2	2
g15	2	2	1	2	2	g230	1	1	1	1	1	g306	1	1	2	1	2
g150	2	2	2	2	2	g233	2	1	1	1	1	g307	1	2	2	2	2
g151	2	2	2	2	2	g234	2	1	1	1	1	g308	1	2	2	2	2
g152	2	2	2	2	2	g237	2	1	1	1	1	g309	1	2	2	2	2
g155	2	2	2	2	2	g239	1	1	1	1	1	g314	1	2	2	2	2
g156	2	2	2	2	2	g241	1	1	1	1	1	g318	1	2	2	2	2
g159	1	1	1	1	2	g242	1	1	1	1	1	g326	1	2	2	2	2
g160	2	2	2	2	2	g244	1	1	1	1	1	g333	2	2	NA	2	2
g162	2	2	2	2	2	g245	1	1	1	1	1	g337	1	2	2	2	2
g17	2	2	1	2	2	g249	1	1	1	1	1	g339	1	2	2	2	2
g170	2	2	1	2	2	g250	1	1	1	1	1	g340	1	2	2	2	2
g173	2	2	2	2	2	g252	2	2	2	2	2	g342	1	2	2	2	2
g176	2	2	2	2	2	g256	1	2	NA	99	2	g343	2	2	2	2	2
g178	2	2	2	2	2	g257	2	1	2	2	2	g38	2	2	2	2	2
g179	1	2	2	2	2	g259	2	2	2	2	2	g40	2	2	2	2	2
g181	2	2	1	2	2	g26	2	2	2	2	2	g44	2	2	2	2	2
g188	2	2	2	2	2	g261	1	2	2	2	2	g48	2	2	1	2	2
g19	1	2	2	2	2	g265	2	2	2	2	2	g7	2	2	1	2	2
g191	2	2	2	2	2	g266	1	2	2	2	2	g72	2	2	1	2	1
g192	2	2	2	2	2	g267	2	2	2	2	2	g73	2	2	1	2	1
g193	2	2	2	2	2	g269	2	2	2	2	2	g74	1	1	2	1	1
g199	1	2	2	2	2	g270	1	2	2	2	2	g75	2	1	1	1	1
g201	2	1	1	1	1	g271	1	2	2	2	2	g77	1	2	1	1	1
g203	2	1	1	1	1	g272	2	2	2	2	2	g8	2	2	2	2	2
g205	2	1	1	1	1	g274	2	2	2	2	2	g80	2	2	1	2	1

**Table 7: MB consensus result. A majority vote value of "99" indicates that the predictions from the three datasets are all different on this subject.**

ID	CA125	MS	NMR	Vote	TRUE	ID	CA125	MS	NMR	Vote	TRUE	ID	CA125	MS	NMR	Vote	TRUE
g111	1	1	1	1	1	g208	1	1	1	1	1	g275	3	NA	2	99	2
g113	2	3	1	99	1	g209	1	1	1	1	1	g276	1	3	3	3	2
g119	1	1	1	1	1	g21	1	2	1	1	3	g277	1	3	2	99	2
g122	1	1	1	1	1	g210	1	1	3	1	1	g28	3	3	3	3	3
g124	1	1	3	1	1	g213	1	1	1	1	1	g281	1	2	2	2	2
g126	1	1	1	1	1	g214	1	NA	1	1	1	g286	1	2	2	2	2
g129	2	1	1	1	1	g215	2	1	1	1	1	g289	3	2	2	2	2
g135	3	3	3	3	3	g216	1	1	1	1	1	g291	1	2	2	2	2
g139	1	3	3	3	3	g217	1	1	1	1	1	g293	1	2	2	2	2
g140	3	1	1	1	3	g220	1	1	1	1	1	g295	1	2	2	2	2
g142	3	1	3	3	3	g221	2	1	1	1	1	g296	1	2	2	2	2
g143	1	1	1	1	3	g223	2	1	1	1	1	g298	1	2	2	2	2
g144	1	2	3	99	3	g225	1	1	1	1	1	g301	1	2	2	2	2
g147	3	3	3	3	3	g227	1	1	1	1	1	g303	1	2	2	2	2
g149	1	3	3	3	3	g23	1	3	3	3	3	g305	2	2	2	2	2
g15	1	3	3	3	3	g230	1	1	1	1	1	g306	1	2	2	2	2
g150	3	3	3	3	3	g233	1	1	1	1	1	g307	1	2	2	2	2
g151	3	3	3	3	3	g234	2	1	3	99	1	g308	1	3	2	99	2
g152	2	2	3	2	3	g237	2	1	1	1	1	g309	1	2	2	2	2
g155	3	3	3	3	3	g239	1	1	2	1	1	g314	1	2	2	2	2
g156	1	3	3	3	3	g241	1	1	1	1	1	g318	1	2	2	2	2
g159	1	1	1	1	3	g242	1	1	1	1	1	g326	1	3	2	99	2
g160	3	3	3	3	3	g244	1	1	1	1	1	g333	2	2	NA	2	2
g162	3	3	1	3	3	g245	1	1	1	1	1	g337	1	2	2	2	2
g17	1	3	1	1	3	g249	1	1	1	1	1	g339	1	1	2	1	2
g170	3	3	1	3	3	g250	1	1	1	1	1	g340	1	2	2	2	2
g173	3	3	3	3	3	g252	1	3	3	3	3	g342	1	2	2	2	2
g176	3	3	3	3	3	g256	1	2	NA	99	2	g343	1	2	2	2	2
g178	3	3	3	3	3	g257	1	2	2	2	2	g38	3	3	2	3	3
g179	1	3	1	1	3	g259	3	3	2	3	2	g40	3	3	3	3	3
g181	3	3	1	3	3	g26	3	3	1	3	3	g44	3	3	1	3	3
g188	3	3	2	3	3	g261	1	2	2	2	2	g48	3	3	1	3	3
g19	1	3	3	3	3	g265	1	2	2	2	2	g7	1	3	3	3	3
g191	3	3	3	3	3	g266	1	2	2	2	2	g72	1	1	1	1	1
g192	3	3	3	3	3	g267	2	3	3	3	2	g73	1	2	1	1	1
g193	2	3	3	3	3	g269	1	2	2	2	2	g74	1	1	3	1	1
g199	1	2	2	2	2	g270	1	3	2	99	2	g75	1	1	1	1	1
g201	2	1	1	1	1	g271	1	2	2	2	2	g77	1	3	1	1	1
g203	2	1	1	1	1	g272	1	2	2	2	2	g8	2	3	1	99	3
g205	1	1	1	1	1	g274	3	1	2	99	2	g80	1	2	1	1	1

**Table 8: 3-Class consensus result. A majority vote value of "99" indicates that the predictions from the three datasets are all different on this subject.**