

Comparisons between Supervised models for cell classification using single-cell RNA-sequence data

Mian Wu¹, Qiushuo Cheng², Yi Shen³, Congjia Chen⁴

1 School of Computer Science, Department of Electronic Information and Electrical Engineering, Dalian University of Technology, Dalian, Liaoning 116024, China

2 School of Marxism, Northeastern University, Shenyang, Liaoning 110169, China

3 School of Mechanical Engineering, Shanghai Jiao Tong University, Shanghai 200240, China

4 School of Life Science, University of Liverpool, Liverpool, Liverpool L69 3BX, United Kingdom

ABSTRACT. *The ability to classify cells based on scRNA-Seq is highly applicable and promising to biology studies. However, the specific computational challenges for biological data induced by high dimension and batch effect may influence the accuracy of cell classification. Based on the previous research that supervised models are more superior in several applications compared with unsupervised models, empirical comparisons between 4 supervised models for cell classification using single-cell RNA-sequence data were finished to evaluate the performance of every supervised model in biology. To overcome the curse of dimensionality, the first step is data pre-processing which includes cell selection, feature selection containing low variance filter and chi square filter, and Principle components analysis. The second step is supervised learning classification using 3 Decision-tree based models and logistic regression. Finally, the comprehensive comparisons of different classification models relied on different standards of evaluation will be done. The results show that logistic regression behaves better than other tree-based models both in accuracy score and f1-score. Our works provide some directions for the supervised model choice in the domain of biology.*

KEYWORDS: *Cell type classification; Machine learning; Computational biology; Supervised learning; Comparison*

1. Introduction

Single-cell RNA-sequencing technology have been rapidly developed in recent years, as a technique which can measure the transcriptome and gene expression level

of individual cell, scRNA-seq can reveal many potential properties of cell subpopulations which could not be accomplished in bulk RNA sequencing [1]. From the count of publications in PubMed (Figure 1), the publication of the scRNA research is increasing dramatically, indicating the remarkable attention worldwide.

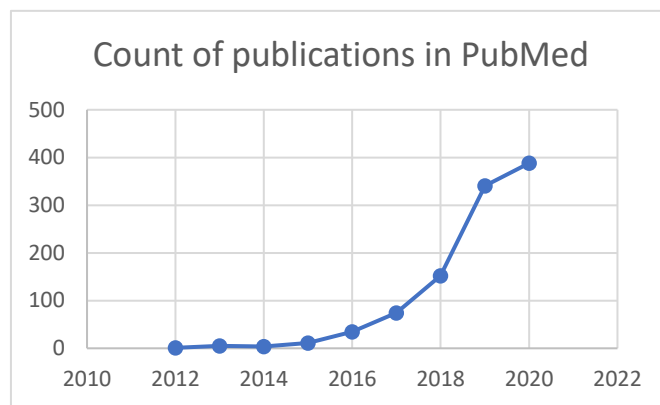


Figure.1 The count of publications about the topic Single-cell RNA-Seq in PubMed

The focus of recent work is on the cell characterization and differentiation within each population being compared. Up to now, the work primarily depended on unsupervised methods or known markers. Known markers, in biological cases, are the specific genes which would be highly expressed in certain types of cells [2]. While the application of markers is useful, it may not be available for several cell types [3]. Although unsupervised methods are useful to solve the analysis of unlabeled dataset, some experts have shown that supervised methods can be much easier to interpret and much more accurate [4]. However, the lack of the empirical evaluation among multiple supervised models applied in single-cell RNA sequencing data will make the supervised learning methods less credible in the biology.

Basically, the analysis of scRNA-seq data is a supervised, high dimensional, multiple classification problems. For the biological gene data, expression level of a single cell type under different experimental situation (termed as batch effect) may be highly variable which is either because of the difference from the sequencing

platform or the variance induced by biological dynamic [5]. Below we will show various methods trying to overcome the problems, try to restore the real situation as much as possible, and solve the overfitting problem. Inspired by a paper published in 2006 [6], which is about the empirical comparison based on 11 binary classification problems, decision trees or the extensions of decision trees are assumed to be the most effective in the prediction of gene expression cell types. However, in terms of the final accuracy score, the optimal result of our project is not from DT, but from logistic regression.

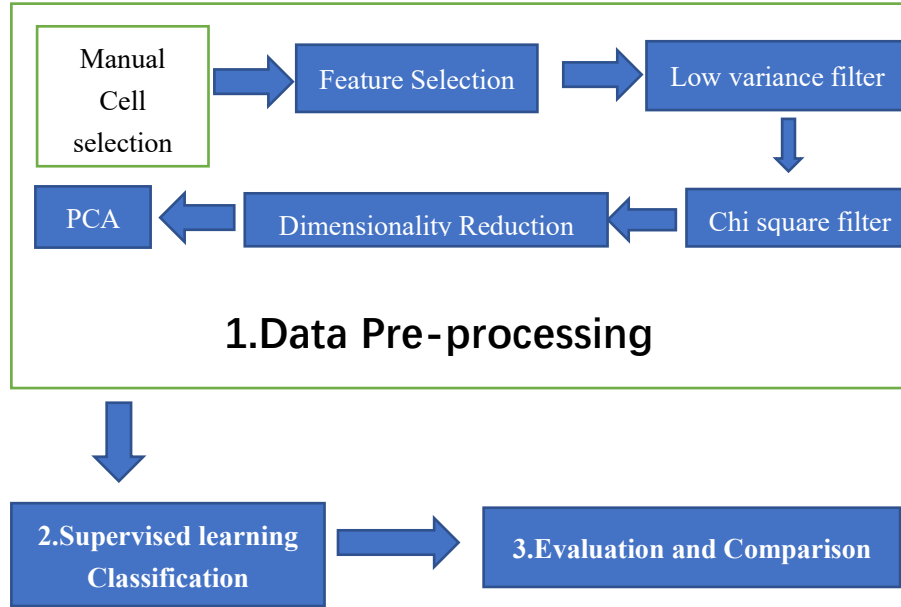
2. Methods

2.1 Data source and pre-processing

Data from mouse gene dataset of interest which contains 20,499 genes with normalized gene expression levels - RPKM values for each cell and 24244 total samples were used [3]. The data contains all_data.h5, together with separate train & test dataset. Classifier for cell types, in real experiment, would be used to predict the cells from new experiments. Therefore, to restore the actual situation, labels that are in the testing dataset can also be found in the training data. Meanwhile, compared with the train samples, test samples come from different experiments were selected.

Table.1 The dataset

	Train dataset	Test dataset
Samples	21389	2855
Unique cell types	46	21
Features (genes)	20499	20499



2.1.1 Cell selection

We select the cell types based on the test cell types manually as our real training labels to minimize the noise given by the train dataset redundancy.

2.1.2 Feature selection

Uninformative feature caused by irrelevance, correlation, and redundancy can impede the performance of classification model. Additionally, because of the technical variation which is mainly caused by the difference of sequencing platform and uninformative biological variation induced mainly by the experiment batch effect, we applied low variance filter based on scikit-learn package [7] to select highly variable genes based on the RPKM value. Meanwhile, chi square function was applied to measure the relation between features and labels.

Table.2 Example of Chi Square table of gene dataset

	Gene 1	Gene 2
Cell 1	RPKM value	RPKM value
Cell 2	RPKM value	RPKM value

Chi Square Formula:

$$\chi_e^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

In the Chi Square Formula, ‘O’ is the observed value, ‘E’ means the expected value and ‘i’ is the ‘i’th position in the contingency table. The chi-squared statistic is a single number that tells you how much difference exists between your observed rpk values and the rpk values you would expect if there were no relationship at all in the dataset. P-value was used here to explain the result of Chi square for different groups. The Select K-Best package was used here to select the top related features [7]. After the feature selection, highly variable genes and highly specific cell type related features would be left. Selection could facilitate downstream applications like DT-based classification and save the computational costs.

2.1.3 Principal component analysis (PCA)

After feature selection, it is still important to apply dimension reduction as too much noise still exists. The core concept of PCA is to map all features to K dimensions. Therefore, PCA for dimensionality reduction was applied and adjusted to the optimal dimension leading to best accuracy.

2.2 Supervised learning methods

2.2.1 Decision Tree

Decision tree is a basic model which was used as a controlled trial in the comparisons. The max-depth and pruning were not determined in our experiment. In our model, we use Gini impurity as our criterion.

2.2.2 Random Forest

In random forests, each DT in the model is built from a set of samples drawn with replacement from the training set, i.e., .632bootstrap, which means roughly 63% of the original data are selected. The input is the entire original training dataset. We use cross validation to determine the best number for trees and the max-depth of the tree in the model. In our model, we use Gini impurity as our criterion.

2.2.3 *Ada-Boosting*

Ada-Boosting uses a set of weak classifier, in our case, small DTs, to operate on repeatedly modified versions of the data. The predictions from all of them are then combined through a weighted majority vote to produce the final prediction. Therefore, there are two sets of weights: weights for DT and weights for data. Initially, those weights are all set to $1/N$. We train the first weak classifier and focusing on the mistakenly classified cells by allocating new weight for each data point. Each subsequent weak learner is thereby forced to concentrate on the examples that are missed by the previous ones in the sequence. Eventually, the training error will reduce to zero.

2.2.4 *Logistic Regression*

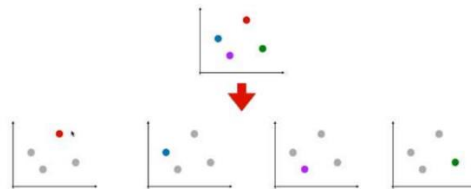


Figure.3 One versus rest

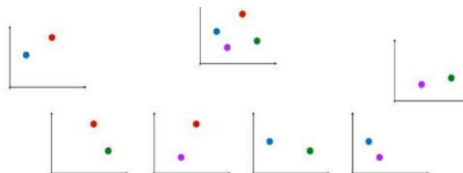


Figure.4 One versus one

We used One vs One (OVO) instead of One vs Rest (OVR) in our logistic regression model. One vs Rest (OVR) treat a multi-class(n) problem as n binary problem. Each binary problem we select one set of data points as one class, and all the other point as the other class. Therefore, the algorithm takes $O(nT)$, if the binary classification takes $O(T)$.

One vs One (OVO) treats a multi-class problem as 2 combinations of n binary problems where N is the number of our classes. For each binary problem we select two sets of data points as two classes. Then we get 2 combinations of n classifiers.

We applied test data to classifiers and make majority votes.

We are using L2 regulation, which means to minimize the function below:

$$\min_{w,c} \frac{1}{2} w^T w + C \sum_{i=1}^n \log(\exp(-y_i(X_i^T w + c)) + 1)$$

2.3 Evaluation of classifications

We used different standards to evaluate our models as they can judge our models from different aspects. We built a confusion matrix and calculated the accuracy score, recall score, precision score and f1_score based on it. Accuracy score is based on the whole data, and the other standards are based on each cell type. In addition, we normalized our confusion matrix to see the accuracy of specific cell types clearly.

However, in evaluation of biology classification models, the question of how similar two cell types are is quite important because a rigid (binary) distinction between cell types is not appropriate since “neuron”, “hippocampus”, and “brain” are all related cell types, and a model that groups these cell types together should not be penalized as much as a model that groups completely unrelated cell types together [3].

Therefore, we improved our evaluation by add some weights. We downloaded a similarity matrix [3]. In general, we changed numbers in the confusion matrix by timing them with the weight which equals to (1-similarity number). Then we acquired our new accuracy score by our weighed confusion matrix.

2.4 Software packages

Low variance filter, Chi square filter, PCA, Decision tree, Random forest, Ada-boosting, Logistic regression, and confusion matrix are all coded in Python based on scikit-learn package [7].

3. Results

3.1 Accuracy score of the four models

3.1.1 Decision Tree

First, let us see how well Decision Tree performed in this problem. In the Decision Tree Model, we both tried model without PCA or with PCA from 40 to 100. Finally, we used 34% as our baseline for comparison between decision tree-based models.

Table.3 Accuracy score of Decision Tree

Chi Square Filter	Lower Variance Filter	PCA	Accuracy score
None	15	100	0.212
None	15	50	0.194
None	15	45	0.262
None	15	40	0.273
None	None	None	0.349

Table.4 Accuracy score of Random Forest

Chi Square Filter	Lower Variance Filter	PCA	Accuracy score
None	15	200	0.321
None	15	100	0.300
None	15	50	0.429
None	15	45	0.430
None	15	40	0.451

3.1.2 *Random Forest*

In the Random Forest Model, the result is better, achieving 45% overall accuracy. The training data for each tree comes from bootstrap. The best result comes from 40-dimension 150 trees with maximum depth of 30.

The hyperparameters are determined using cross validation, for example, to determine the number of trees in the random forest, we perform 5- fold cross validation on 100 trees, 150 trees 200 trees and so on. The high cross validation mean accuracy comes from 150 trees, which is 93.2%

3.1.3 *Ada-Boosting*

Table.5 Accuracy score of Ada-Boosting

Dimension after PCA	N_estimator	Max_depth	Accuracy on training set	Accuracy on test set
50	100	3	0.30	0.18
50	300	3	0.42	0.18
50	1000	3	0.47	0.21
50	100	10	0.89	0.26
50	300	10	0.94	0.32
50	800	10	0.95	0.33
40	300	10	0.94	0.38
40	400	10	0.95	0.39
40	850	7	0.85	0.37
40	1400	7	0.90	0.39

We selected some of the most representative parameter values that we tried to show you in Table 5. From the last two rows in the table, you can see that even if we almost double the max iteration number(n_estimator), the overall training error improve only slightly. It becomes incredibly time-consuming. By adjusting the Decision Tree's max depth, we manage to improve the result, but still, the highest success rate cannot even compete with random forest classifier's worst result.

Therefore, our conclusion is that DT-based boosting algorithm is not suitable to be directly applied in a multi-class biology classification problem.

3.1.4 Logistic Regression

Table.6 Accuracy score of Logistic Regression

Solver	Multiclass	Max_iter	Accuracy on training set	Accuracy on test set
Before Cell Selection				
Sag	Multiclass	100	0.78	0.5
Sag	Multiclass	200	0.81	0.54
Sag	Multiclass	500	0.84	0.54
Sag	Multiclass	1000	0.86	0.52
Sag	Multiclass	2000	0.88	0.49
After Cell Selection				
Sag	Multiclass	100	0.88	0.55
Sag	Multiclass	200	0.9	0.57
Sag	Multiclass	500	0.9	0.56
Sag	Multiclass	1000	0.93	0.55
Sag	Multiclass	2000	0.93	0.51

Logistic Regression works best. It is the most time-efficient and has the best overall accuracy. As described in the method, the parameters chosen are for One vs One (OVO) multiclass classification, which is more time consuming than One vs Rest (OVR), but at the same time more accurate. If we add the cell selection, the result is even better, achieving 57% success rate.

3.2 Confusion matrix on random forest and logistic regression

Besides accuracy score, we also used confusion matrix, precision, recall, and f1-score (Methods) to evaluate our models. We mainly compared the two best models in terms of overall accuracy. In the figure 5 and figure 6, the deeper the color of each intersection grid is, the more cases there are that X label is predicted to be Y

label. From the figure 5 and figure 6, we can see that there are more deep color grids in the diagonal in the Logistic Random than in the Random Forest, which means that more labels in the Logistic Regression were predicted correctly.

In fact, we can see that the number of cell-types each time the confusion matrix presents depends on the union of cells from both test data, which has 21 types of cells, and the predicted cell dataset. In that case, the confusion matrix, which is output from each model, or from the same model but at different times, will vary in the number of label-types. However, that does not affect our conclusion one thing. That is because the only correct case of classification is when a labeled cell is predicted to be itself, which is denoted by the diagonal grids, and those deep color grids in the diagonal only come from the 21 cell types to be predicted. Therefore, more labels to be predicted correctly means more labels to be predicted correctly within those 21 types of cells.

As well, we can also see that there are more deep color grids in the lower left corner of Random Forest, which means that Random Forest assigns more labels to the wrong kinds.

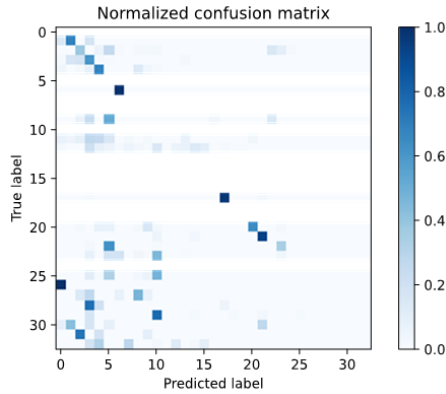


Figure.5 Random Forest

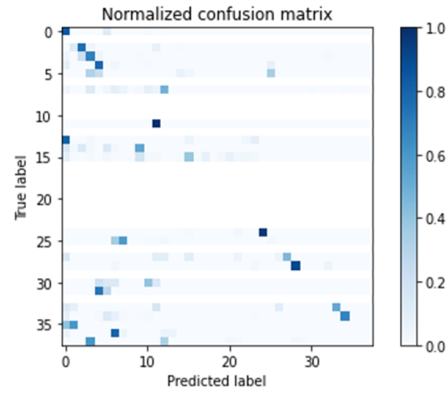


Figure.6 Logistic Regression

3.3 Precision, recall, f1-score on random forest, and logistic regression

Among all the test cells, we test the precision score, recall score and f1-score for them. We also provided the number of each type of cells in the test data, so that we can have a more intuitive insight into the relationship between number of cells and its prediction result. We have already known that f1-score is the harmonic average value of precision score and recall score (Methods).

Table.7 Random Forest

	precision	recall	f1-score	support
CL:0000137 osteocyte	0.98	0.97	0.98	108
CL:0000235 macrophage	0.85	0.93	0.89	42
UBERON:0000966 retina	0.99	0.65	0.79	250
CL:0002319 neural cell	0.58	1	0.74	81
UBERON:0001003 skin epidermis	0.83	0.63	0.72	678
CL:0002321 embryonic cell	0.6	0.69	0.64	173
CL:0002322 embryonic stem cell	0.5	0.6	0.55	358
UBERON:0000044 dorsal root ganglion	0.3	0.37	0.33	123
CL:0000037 hematopoietic stem cell	0.59	0.06	0.11	162
UBERON:0001851 cortex	0.16	0.04	0.07	266
UBERON:0001264 pancreas	1	0.01	0.02	162
CL:0000353 blastoderm cell	0	0	0	10
CL:0000540 neuron	0	0	0	133
CL:0000746 cardiac muscle cell	0	0	0	11
UBERON:0000115 lung epithelium	0	0	0	78
UBERON:0000922 embryo	0	0	0	60
UBERON:0000955 brain	0	0	0	38
UBERON:0001898 hypothalamus	0	0	0	29
UBERON:0001954 Ammon's horn	0	0	0	15
UBERON:0002048 lung	0	0	0	58
UBERON:0002107 liver	0	0	0	20

accuracy			0.43	2855
----------	--	--	------	------

Table.8 Logistic Regression

	precision	recall	f1-score	support
CL:0000137 osteocyte	1	0.98	0.99	108
CL:0000235 macrophage	0.89	0.98	0.93	42
UBERON:0001003 skin epidermis	0.79	0.89	0.83	678
UBERON:0000955 brain	0.96	0.71	0.82	38
UBERON:0000966 retina	0.96	0.69	0.8	250
UBERON:0002107 liver	0.93	0.7	0.8	20
CL:0000037 hematopoietic ste m cell	0.66	0.68	0.67	162
CL:0002321 embryonic cell	0.68	0.64	0.66	173
CL:0002319 neural cell	0.41	1	0.58	81
UBERON:0000044 dorsal root ganglion	0.47	0.76	0.58	123
CL:0002322 embryonic stem cell	0.43	0.62	0.51	358
UBERON:0001851 cortex	0.22	0.15	0.18	266
CL:0000540 neuron	0.21	0.06	0.09	133
UBERON:0001264 pancreas	1	0.01	0.02	162
accuracy			0.57	2855

未加入细胞筛选

	precision	recall	f1-score	support
CL:0000137 osteocyte	1	0.96	0.98	108
CL:0000235 macrophage	0.9	0.9	0.9	42

UBERON:0001003 skin epidermis	0.79	0.87	0.83	678
CL:0002319 neural cell	0.68	1	0.81	81
UBERON:0000966 retina	0.99	0.67	0.8	250
UBERON:0000044 dorsal root ganglion	0.74	0.79	0.76	123
UBERON:0000955 brain	1	0.53	0.69	38
CL:0002321 embryonic cell	0.6	0.76	0.67	173
CL:0002322 embryonic stem cell	0.62	0.67	0.64	358
UBERON:0002107 liver	1	0.45	0.62	20
CL:0000037 hematopoietic stem cell	0.86	0.38	0.53	162
UBERON:0001851 cortex	0.13	0.05	0.07	266
CL:0000540 neuron	0.14	0.02	0.04	133
CL:0000353 blastoderm cell	0	0	0	10
CL:0000746 cardiac muscle cell	0	0	0	11
UBERON:0000115 lung epithelium	0	0	0	78
UBERON:0000922 embryo	0	0	0	60
UBERON:0001264 pancreas	0	0	0	162
UBERON:0001898 hypothalamus	0	0	0	29
UBERON:0001954 Ammon's horn	0	0	0	15
UBERON:0002048 lung	0	0	0	58
accuracy			0.54	2855

We mainly compared Random Forest and Logistic Regression. As can be seen from Table 7 and Table 8, all values that are greater than 0.9 are bolded. From the f1-score column, we can see that, osteocyte cells were always classified the most accurately, and macrophage cells the second. We can also notice that, f1-score of Logistic Regression is usually larger than Random Forest, which means that Logistic

Regression has a better classification ability. In conclusion, Logistic regression behaves better than random forest.

3.4 How pre-processing impact our results

Both feature selection and dimensionality reduction tools showed great impact on our models. We compared different combinations of these methods to pre-process our data. In the Decision Tree Model, we found that Decision Tree without PCA (Principal Component Analysis) performed much better than with PCA (Principal Component Analysis). Presumably, the model was relatively simple, therefore, it will be unable to effectively differentiate between cells when there were relatively fewer features for it to learn. In the Ada-Boosting Model, PCA (Principal Component Analysis) seemed to have very little impact on the accuracy score. However, feature selection had a very large impact on the program efficiency. Considering its low classification accuracy results, we mainly compared the best two classifiers in our experiment- Random Forest and Logistic Regression. Compared with the Logistic Regression Model, which though has the best accuracy score, we found that pre-processing had a larger impact on our Random Forest Model, especially when combined, which improved the accuracy score of Random Forest by 10% on average, Logistic Regression by 5% (See Table 9). From above, we can see that when our models were getting relatively more complicated compared with Decision Tree, the pre-processing could have a larger impact on the accuracy of the models.

Table.9 Pre-processing impact on the accuracy score

Feature selection types	Accuracy increase for RF	Accuracy increase for LR
Chi square filter (PCA 40)	2%	2%
Low variance filter (PCA 40)	5%	3%
Two-combined (PCA 40)	10%	5%

3.5 How similarity between cells revises our results

We used the similarity coefficients to revise the accuracy of the two best models. For the Random Forest Classifier, the accuracy score before applying similarity

coefficients was 0.435, and after was 0.461, which increased by 2.6%. For Logistic Regression Model (the best model), the accuracy score before was 0.571, and after was 0.598, which increased by 2.7%. (See Table 10)

In conclusion, we revised the accuracy score according to Similarity Coefficients.

Table.10 Revision of the accuracy score using similarity coefficients between cells

Classifier	Accuracy before	Accuracy after	Accuracy increase
Random Forest	0.435	0.461	2.6%
Logistic Regression	0.571	0.598	2.7%

4. Discussion

Back to our problem: How to improve the accuracy of the high dimension, multiclass, batch-related classification problem? We took what we learn from class, mostly ensemble methods, to implement on our model and I think it is safe to say we overcome this problem in some way, improving from baseline Decision Tree 34.9% to Logistic Regression 57.1% before revision, and 59.8% after revision.

In all the models that we have tried, the Logistic Regression Model performed best, and the Random forest one the second. The Ada- Boosting model seemed to be the most time consuming and having the lowest accuracy. One possible reason why the Logistic Regression outperformed all the other Decision Tree-based models is that Logistic Regression uses a modified version of "divide and conquer". It breaks down multiclass problem into several binary classification problems, and then combines the result. This approach makes the algorithm more time efficient, and the voting technique also improves the Logistic Regression Model.

At the pre-processing stage, we first applied cell type selection by extracting some of the cells which are selected based on the test cell dataset from train dataset, in order to reduce the computational difficulty and increase the accuracy of our supervised models. For the sake of the simplicity and manually filtering noise, we raise the accuracy of our classifiers for 10% on average. However, we are aware of that use of manually selection to reduce the train data set will not fit the real-world experiment. Therefore, we are much more confident to classify the cell type in the

situation that the researchers have the basic prediction and range of the potential cells. We also used two kinds of Feature Selection tools called Chi Square Filter and lower Variance Filter, and PCA (Methods). They all showed great impact on our models, not only increased the accuracy score, but also improved the efficiency of models to a large extent, especially for large-scale scRNA-seq datasets whose computational time is long and memory-consuming. Furthermore, when we were testing the performances of each model with each PCA dimensionality, there was usually a 3% fluctuation, so we usually took the average of the three tests as the final accuracy score.

With regard to Decision Tree Model, as can be seen in table3, we found that Decision Tree without dimensionality reduction performed much better. This is an interesting finding. One possible reason is that the PCA method does not use labels, so the purpose of PCA is for reconstruction rather than classification.

5. Future work

For the best model in our research, the accuracy raised by the feature selection is not dramatic, therefore, other feature selection method such as the selection based on highly expressed gene should be tried in the future. Also, we can also take the advantage of the idea of divide-and-conquer to design an algorithm combined with Ada-Boost and One vs Rest (OVR), instead of just using the existing package, converting the multiclass classification problem into several binary classification problem, to check if the accuracy can be further improved. In addition, we can take into consideration to build a working software or a website server, therefore, it is necessary to consider designing adaptive hyperparameters. For example, the number of trees in random forest can be adaptively updated to the optimal value after new data is added to the training set. Furthermore, we need to generalize our model by discarding cell selection, because in real world, we usually cannot guarantee to extract labels every time we have a new test set. Since our project showed that cell selection did work, in the future, instead of cell selection, we can try some other methods to detect irrelevant cells in the training data to improve the result.

Conclusion

References

- [1] ALQUICIRA-HERNANDEZ, J., SATHE, A., JI, H. P., NGUYEN, Q. & POWELL, J. E. 2019. scPred: accurate supervised method for cell-type classification from single-cell RNA-seq data. *Genome Biology*, 20, 264.
- [2] EL AMRANI, K., ALANIS-LOBATO, G., MAH, N., KURTZ, A. & ANDRADE-NAVARRO, M. A. 2019. Detection of condition-specific marker genes from RNA-seq data with MGFR. *PeerJ*, 7, e6970-e6970.
- [3] ALAVI, A., RUFFALO, M., PARVANGADA, A., HUANG, Z. & BAR-JOSEPH, Z. 2018. A web server for comparative analysis of single-cell RNA-seq data. *Nature Communications*, 9, 4768.
- [4] LIN, C., JAIN, S., KIM, H. & BAR-JOSEPH, Z. 2017. Using neural networks for reducing the dimensions of single-cell RNA-Seq data. *Nucleic Acids Research*, 45, e156-e156.
- [5] TRAN, H. T. N., ANG, K. S., CHEVRIER, M., ZHANG, X., LEE, N. Y. S., GOH, M. & CHEN, J. 2020. A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biology*, 21, 12.
- [6] CARUANA, R. & NICULESCU-MIZIL, A. 2006. An empirical comparison of supervised learning algorithms. *Proceedings of the 23rd international conference on Machine learning*. Pittsburgh, Pennsylvania, USA: Association for Computing Machinery.
- [7] PEDREGOSA, F., VAROQUAUX, G., GRAMFORT, A., MICHEL, V., THIRION, B., GRISEL, O., BLONDEL, M., PRETTENHOFER, P., WEISS, R., DUBOURG, V., VANDERPLAS, J., PASSOS, A., COURNAPEAU, D., BRUCHER, M., PERROT, M. & DUCHESNAY, É. 2011. Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.*, 12, 2825–2830.