

1 1.1 Exploratory data analysis

1.1 Basic facts

We have a dataset that consists of 5471 (n) samples with 4124 (p) columns. The features of our data are all continuous, in a logarithmic scale and they expression levels for genes. It is worth mentioning that the dataset is sparse, that is, a lot of cells have many gene expressions that are 0. While it is hard to have an objective measure of sparsity we can plot average gene expression levels across all our features, see Figure 1 below, most genes are non zero for a limited number of cells.

Figure 1: Possible figure of average gene expression level across all genes. We can do either histogram or histogram by label.

We have 2 distinct classes of cells, the TREG cells and the CD4+T cells. These are going to be our labels that we want to classify. We have some class imbalance, the ratio is 6/10 in favour of the CD4+T which is the dominant class. This is going to be important for the models we try to tune as some models have options to adjust for class imbalance. We will also try the option of tuning the threshold decision for classifying to one cell or another. By default the probabilistic models in Scikit-learn classify to the positive class if the conditional probability for the given model $\mathbb{P}(y|X) > 0.5$ ¹. We will tinker with this threshold to try to optimize the F1-Score given our class imbalance.

1.2 Visualization and dimensionality reduction

Both the fact that $n \approx p$ and the fact that the data is sparse point us towards using regularization and feature selection. The instructions for this problem also make us use PCA with 10 components. To confirm whether the number of components is optimal, we can plot the cumulative sum of the explained variance by each of the components. As you can see below in Figure 2, choosing just the first ten components makes us use a very amount of the total variance. In Section 3 we will try to tune the number of components to get an improved F1-Score.

Figure 2: Scree plot: cumulative sum of PCA components.

For completeness, we add Figure 3 where we use t-SNE to reduce the dimensionality of the data from 4124 to 2^2 with the purpose of visualizing the joint distribution (after the t-SNE transformation) of both cells. Figure 3 shows some separability of the two classes. It also shows some cells of a given class (TREG) in regions where the density is much higher for the other type of cell (CD4T), maybe we can use this as intuition as to why in later sections we find it difficult to improve the F1-Score beyond 0.95.

Figure 3: Joint distribution of transformed data by t-SNE

2 1.2

- Most tuning is done, VSM give very good results.
- Tough to do what shak is doing with pca Example from sklearn
- RESULTS FROM GRIDSEARCH ARE NOT GOING TO BE REPLICABLE, RESULTS FROM FINAL ESTIMATOR BASED ON IT SHOULD BE.

3 Our 3 models

- AdaBoost, just because we have seen it in class and it could be interesting.

¹see scikit references

²The underlying algorithm is stochastic and quite sensitive to how we tune the hyperparameter of perplexity

- I wanted some classically strong classifier like random forest or GBDT?
- Play around with the PCA optimality. Or some other modern dimensionality reduction tool like t-SNE or UMAP.

4