

Statistical Significance of Cluster Membership

Neo Christopher Chung*

January 2018

Contents

Abstract	2
Introduction	3
Methods and Algorithms	4
Results	7
Simulation Studies	7
Genomic Applications	8
Discussion	9
Figures	10
Supplementary Figures	14
References	17

*Institute of Informatics, Faculty of Mathematics, Informatics and Mechanics, University of Warsaw.
Contact Information: <http://ncc.names>

Abstract

Clustering is routinely applied to modern high-dimensional data, including gene expression measurements from microarray and RNA-seq. Iteratively estimating the cluster centers and assigning memberships according to pre-defined criteria, the clustering algorithms classify genes or samples to help ascertain molecular processes or sub-types. For example, the cluster membership assignments of unlabeled single cells from massively parallel RNA-seq experiments are used as the cell identities. However, how can we evaluate if the cluster memberships are correctly assigned? To this end, we introduce the jackstraw methods for unsupervised classifications that rigorously test the assignments of data features into their clusters. By learning uncertainty in clustering the noisy data, the proposed jackstraw methods can identify statistically significant features that truly make up the corresponding clusters. Simulation studies using K -means clustering confirm the accuracy of the proposed statistical significance. We consider mRNA abundances of 5981 *Saccharomyces cerevisiae* genes under cell cycle. After the proposed jackstraw methods are applied for $K = 6$ clusters, we estimate and use posterior inclusion probabilities (PIP) to select and visualize the canonical features for their clusters. We also investigate the single cell RNA-seq (scRNA-seq) data from a mixture of Jurkat and 293T cell lines, where individual cell identities are unknown. The jackstraw methods evaluate cluster membership assignments of 3381 unlabeled single cells such that the majority of multiplets are identified in an unsupervised manner. When clustering is employed in high-dimensional data analysis, the proposed tests enable rigorous evaluation of membership assignments that readily improve feature selection and visualization.

Software: jackstraw package in R available at <https://github.com/ncchung/jackstraw>.

Abbreviations: single cell RNA sequencing (scRNA-seq), principal component analysis (PCA), posterior inclusion probability (PIP), false discovery rate (FDR)

Introduction

High-throughput technologies have enabled large-scale measurements of DNA, RNA, metabolites, and others. Recent technological and experimental advancements, such as single cell RNA-seq [1, 2] and mass-spectrometry [3, 4], have resulted in increasing challenges and opportunities for using unlabeled data and unsupervised learning. For example, a single cell RNA-seq (scRNA-seq) technology enables gene expression measurements of thousands of blood cells in order to elucidate molecular subtypes. Unsupervised assignments of unlabeled single cells to K clusters according to their gene expression profiles provide cluster-based cell identities. Despite diverse clustering techniques available, it has not been possible to re-use data-driven clusters and test their membership in downstream statistical analyses without incurring artificially inflated significance. We have developed a novel and general method to statistically test the assignment of a data feature to a particular cluster, aiding in feature selection, dimension reduction, and visualization.

Clustering has been one of the most popular analysis methods for high-dimensional genomic data. In the absence of external and accurate labels, clustering can identify and approximate co-regulated subsets of genomic variables (e.g., genes, loci) or subtypes of related observations (e.g., patients, single cells). For example, a conventional microarray study measures gene expression of samples from either control or disease groups. Since the molecular functions of genes might not be known, the unsupervised classification of genomic variables can help identify co-varying subsets that form molecular processes [5–7]. These clusters and membership assignments of genes have been extensively used in the visualization of systematic patterns and outliers [8, 9]. Recently, there have been many studies where mRNA abundances from thousands of single cells are measured en masse using scRNA-seq [10–12]. Then, gene expression profiles of unlabeled single cells are clustered to obtain cell identities. These cell identities, which may be related to subtypes, lineage, or other molecular factors, are often used in downstream differential expression and other analyses. Note that a data feature refers to either a variable or an observation, since clustering can be applied on either dimension.

After automatically assigning observed features to K clusters that are summarized by K centers, we are interested in testing the membership assignments of individual features. This will improve the data-driven cell identities in scRNA-seq experiments, as well as the clustering of genomic variables to help elucidate molecular processes. To this end, we have developed an innovative data resampling and testing scheme for unsupervised classification that rigorously evaluates whether observed features are truly members of corresponding clusters. By estimating p-values and posterior inclusion probabilities (PIPs), the proposed methods can identify and visualize features that have been accurately and reliably assigned to the clusters. This bridges direct estimation of latent variables from large-scale data and fundamental hypothesis framework, which readily provides p-values, false discovery rates, and posterior probabilities crucial for data exploration and inference.

By utilizing a newly developed resampling technique called the jackstraw [13], the proposed methods learn overfitting inherent in using cluster centers that are estimated from the data. In other words, the proposed methods enable the accurate statistical testing of cluster membership while taking into account uncertainty in the clustering algorithms. Simulation studies demonstrate accurate and favorable operating characteristics. The joint behavior of p-values are scrutinized by conducting 100 independent simulations that satisfy the joint null criterion [14]. Two applications are presented using two different dimensions (genomic variables and samples) available for unsupervised classification. Yeast cell cycle microarray data [8] are used to cluster 5981 genes into $K = 6$ clusters, whose statistically significant members are identified. We also consider the scRNA-seq data from a mixture of two different cell lines [12]. By applying the proposed methods on unlabeled single cells, we show improved classification and visualization of cell identities. These proposed methods are implemented in a R package called *jackstraw* (<https://github.com/ncchung/jackstraw>).

Methods and Algorithms

The observed data $\mathbf{Y}_{(m,n)}$ contains m rows and n columns. Because either a set of variables (e.g., genes) or a set of observations (e.g., cells) may be clustered, we refer to m rows as m observed features for simplicity¹. Then, m_1, \dots, m_K features are assigned into corresponding $1, \dots, K$ clusters, where $\sum_{k=1}^K (m_k) = m$. The center \mathbf{c}_k for $k = 1, \dots, K$ summarizes that k^{th} cluster. For example, in K -means clustering, the nearest means are used to assign observed features to the clusters. If \mathbf{y}_i is assigned to k^{th} cluster with \mathbf{c}_k , its membership indicator $\beta_{i,k}$ is 1. By definition, the subset of features \mathbf{y}_i with $\beta_{i,k} = 1$ make up \mathbf{c}_k .

Cluster centers and membership assignments may be viewed as approximating latent variables \mathbf{L} and membership indicators \mathbf{B} (i.e., dichotomous coefficients). Latent variables \mathbf{l}_k for $k = 1, \dots, K$ may assume a wide range of patterns including continuous or categorical structures [15, 16]. Clustering algorithms simultaneously identify the data features that contribute to the estimates of \mathbf{L}_k :

$$\mathbf{Y}_{(m,n)} = \mathbf{B}_{(m,K)} \mathbf{L}_{(K,n)} + \mathbf{E}_{(m,n)}$$

If a particular i^{th} feature is truly associated with a k^{th} latent variable, its coefficient $b_{i,k}$ is 1. Otherwise, 0. Feature-specific noise \mathbf{e}_i is defined as identically and independently distributed. Row-wise means are handled by centering the data, whereas row-wise variances are preserved by our proposed resampling scheme.

There have been important developments in clustering that consider mixture or latent variable models that improve our understanding and interpretation of data [17–19]. However, even model-based clustering approaches or regularization do not provide cluster centers and membership assignments that can be used again against the observed features, resulting in so-called “double dipping.” Our proposed approach learns and incorporates inevitable uncertainty in assigning features to clusters, that are directly derived from the same set of features. This mirrors the jackstraw test when latent variables are estimated using principal component analysis (PCA) [13]. Our statistical significance approach using the jackstraw strategy is related to [20, 21], as well as Bayesian p-values [22, 23]. Furthermore, regularized methods are available for clustering, such that sparsity can be induced [24, 25].

Jackstraw Data and Strategy

We apply the jackstraw strategy to clustering unlabeled features of observed data \mathbf{Y} . Generally, we would like to create a relatively small number s ($\ll m$ or n) of synthetic null features without disturbing the overall patterns of systematic variation. The jackstraw data \mathbf{Y}^* refers to this revised data, where $m - s$ observed features are intact and s synthetic null features have been resampled with replacement (Figure 1). Applying the clustering algorithm on \mathbf{Y}^* produces cluster centers \mathbf{c}_k^* that are almost identical to the original cluster centers \mathbf{c}_k (for $k = 1, \dots, K$).

Because of the nature of clustering algorithms, all features in \mathbf{Y}^* , including s synthetic nulls, will be assigned to one of K clusters. When a synthetic null feature \mathbf{y}_i^* is assigned to k^{th} cluster, an association statistics between \mathbf{y}_i^* and \mathbf{c}_k^* is under the null model that assumes independence since \mathbf{y}_i^* is i.i.d. by definition. Yet, because \mathbf{y}_i^* does indeed contribute to \mathbf{c}_k^* , we effectively learn the overfitting characteristics of the clustering algorithms. Over a large number of iterations $b = 1, \dots, B$, we can form the empirical distribution of null statistics as in *Algorithm 1*.

Feature-level evaluation of cluster membership requires a pre-defined number of clusters K . There is a vast amount of literature on the choice of K , which is beyond the scope of this study. In practice, a data analyst must explore the observed data, often utilizing prior knowledge, visualization, and heuristics. Methods have been proposed in the last five decades in this area of research including cluster stability or reliability statistics [26–34]. We recognize that data normalization, cluster stability, and other pre-classification steps are essential to sensible unsupervised learning. Through re-analysis of microarray and scRNA-seq data, we showcase the jackstraw tests in a context of broader unsupervised learning pipelines.

¹This convention is also followed in the software package where the rows of input data are clustered and tested.

Algorithm 1: Jackstraw Strategy for Unsupervised Classification

- 1 Apply the clustering algorithm to \mathbf{Y} , to obtain \mathbf{C} and β
 - 2 Compute the observed statistics, relating \mathbf{Y} and \mathbf{C}
 - 3 Create \mathbf{Y}^* with a small number of synthetic null features \mathbf{y}^*
 - 4 Apply the clustering algorithm to \mathbf{Y}^* , to obtain \mathbf{C}^* and β^*
 - 5 Compute the null statistics, relating \mathbf{y}^* and \mathbf{C}^*
 - 6 Repeat the above three steps to form an empirical distribution of null statistics
-

There are idiosyncratic outcomes of clustering that require our attention. Some clustering algorithms may generate an empty cluster or a singleton (a cluster with one feature). An empty cluster can be ignored in our methods as it does not contain any observed feature as a member. We consider the only feature of a singleton as its true member. It is possible that synthetic null features are rarely clustered into a certain cluster, such that there is a limited amount of empirical null statistics for that cluster. This likely occurs when that cluster is substantially smaller than others or has very distinct centers such that its members are tightly (and accurately) grouped in n dimensions. An increase of B would alleviate this, in tandem with examining the overall p-value distribution.

Jackstraw Tests for K -means Clustering

We now present a detailed algorithm using K -means clustering [35–37]. K -means clustering is one of the most popular and well-studied algorithms that has been applied to a wide range of genomic studies. In this *Algorithm 2*, we use F-statistics where the full models include appropriate cluster centers. The use of F-statistics allows us to flexibly specify the full and null models, which may incorporate other covariates in more complex settings.

Algorithm 2: Jackstraw Test for Membership Assignments in K -means Clustering

- 1 Apply K -means clustering to the observed data \mathbf{Y} , resulting in cluster centers \mathbf{c}_k for $k = 1, \dots, K$ and membership assignments $b_{i,K}$ for $i = 1, \dots, m$ and $K = 1, \dots, k$
 - 2 Compute the observed statistics F_1, \dots, F_m , where the full models include corresponding cluster centers \mathbf{c}_k
 - 3 Create s synthetic nulls by resampling a small proportion of features $s \ll m$ with replacement, resulting in a jackstraw data \mathbf{Y}^* , with $m - s$ observed features and s synthetic features
 - 4 Apply the clustering algorithm to the jackstraw data \mathbf{Y}^* , resulting in cluster centers \mathbf{c}_k^* and membership assignments $b_{i,K}^*$
 - 5 Compute the null statistics F_1^*, \dots, F_s^* , where the full models include corresponding cluster centers \mathbf{c}_k^*
 - 6 Repeat the above three steps $b = 1, \dots, B$ times to obtain a total $s * B$ of null statistics
 - 7 Compute the p-values by empirically ranking the observed statistics among the null statistics, stratified by cluster assignments
-

The choices of s and B controls the speed of computation, while the total number of null statistics ($s \times B$) determines the overall p-value resolution. For B iterations we need to cluster the jackstraw data B times, and for each iteration $b = 1, \dots, B$, we can obtain s null statistics. Assuming $s \times B$ is hold constant,

a smaller s provides more accurate p-values, while increasing computational burdens. Therefore, we want to ensure the original clusters are preserved as much as possible, permitting the computational power. As we increase the number of synthetic null features s in \mathbf{Y}^* , the overall systematic variation captured by K cluster centers may be substantially disrupted (seen as an increasing proportion of \mathbf{y}^* in Figure 1). While we recommend $s < .1 \times m$ for genomic data, although the number of clusters (K) and the proportion of features assigned to them (m_1, \dots, m_k) must be considered. A higher value of K for a given m would need a smaller s , so that the clusters with limited members are represented in the jackstraw data.

The overwhelming disruption would further inflate null F-statistics, since a larger number of synthetic null features would make up \mathbf{c}_k^* . In extreme scenarios where all features have been resampled, the new cluster centers are completely dominated by independent synthetic null features. This operating characteristic allows us to guard against artificially inflated significance and to guide the input parameters for the proposed algorithm. In practice, we input \mathbf{C} as the initial centers for K clusters when clustering the jackstraw data for efficient convergence. Furthermore, when a computational cost is a concern, one may correlate \mathbf{C} and \mathbf{C}^* to ensure comparability.

In contrast, the conventional resampling methods can be applied to the cluster centers, resulting in a “naive” significance test. After all m features are resampled with replacement, their F-statistics with respect to \mathbf{c}_k are used to form an empirical distribution of null statistics. Observed F-statistics are compared to this empirical distribution to obtain naive p-values. This circular analysis inflates statistical significance, since the observed features are used twice to compute the cluster centers and to again test against the cluster centers. Essentially, this represents how the bootstrap or the permutation approaches would be applied to cluster membership assignments. We apply the conventional methods in simulation studies to demonstrate how the jackstraw approach overcomes this type of circular analysis.

Posterior Inclusion Probabilities

After the membership assignments for k^{th} cluster are tested using the jackstraw, we investigated how to harness their m_k p-values (or, the distribution of null statistics) to filter, de-noise, and visualize the clusters. When considering high-dimensional features typical in large-scale genomic studies, it is advantageous to consider a family of multiple hypotheses simultaneously [38]. Particularly, from m_k jackstraw p-values, we propose to calculate posterior probabilities that features are included in a given cluster. A discussion of posterior inclusion probabilities (PIPs) that are used for shrinkage and improvement of latent variable estimates is available in *Chapter 3* of [39].

Consider that the m_k jackstraw p-values $\mathbf{p}_k = p_{1,k}, \dots, p_{m_k,k}$ are obtained for m_k features that have been assigned to k^{th} cluster. We are interested in estimating a posterior probability that $b_{i,k} \neq 0$, since non-zero coefficients imply their bona fide inclusion in the cluster:

$$\rho_{ik} = \Pr(b_{ik} \neq 0 | \mathbf{p}_{m_k}) = 1 - \Pr(b_{ik} = 0 | \mathbf{p}_{m_k}).$$

PIP can be readily obtained by estimating $\Pr(b_{ik} = 0 | \mathbf{p}_{m_k})$ through an empirical Bayes approach [40,41]. In multiple hypothesis testing, $\Pr(b_{ik} = 0 | \mathbf{p}_{m_k})$ is called a local false discovery rate (FDR). There also exist related Bayesian methods that could be explored for specific applications and prior knowledge [42–44]. These results in m PIPs for K families of multiple hypothesis tests corresponding to K clusters, that can be used for:

1. Retaining a subset of features \mathbf{y}_i with $\rho_i > \alpha_\rho$, where α_ρ is a user-defined threshold,
2. Visualizing features in reduced dimensions (e.g., PCA, t-SNE) where transparency $\sim \rho_i$,
3. Improving the cluster centers by weighting the corresponding features with ρ_i .

Local FDRs and PIPs from K families of multiple hypothesis tests can be flexibly combined for downstream analyses, as to aid feature selection and dimension reduction. When applying the proposed methods on microarray and scRNA-seq data, we incorporate PIPs to hard-threshold and soft-threshold the observed features. Furthermore, this approach may improve a wide range of clustering, by providing probabilistic measures and/or translating into fuzzy clustering algorithms.

Results

Unsupervised classification allows us to non-parametrically cluster large-scale data in absence of accurate external labels for data features. Given the set of features are assigned into K clusters, the proposed methods test their cluster membership assignments. To demonstrate its operating characteristics, we conducted comprehensive simulation studies, which enabled a critical assessment using the underlying truth (*Oracle Groups*). We then applied the proposed methods on a microarray study of *Saccharomyces cerevisiae* that examines the cell cycle and another scRNA-seq data from a mixture of Jurkat and 293T cell lines whose cell identities are of interest.

Simulation Studies

In the simulation studies, we follow the latent variable model described in *Methods and Algorithms*. Latent variables \mathbf{L} are drawn from the Normal($\mu = 0, \sigma^2 = 1$) distribution. Relationships between \mathbf{I}_k and features are given by dichotomous coefficients \mathbf{B} where $b_{i,k}$ indicates whether \mathbf{y}_i is a member of \mathbf{I}_k for $k = 1, \dots, K$ and $i = 1, \dots, m$. The noise \mathbf{B} is drawn i.i.d. from Normal($0, \sigma_b^2$), where its variance governs the noise level. A total of $m = 1000$ features (rows) are simulated over $n = 100$ dimensions (columns). Forming *Oracle Group A*, 500 rows are true members of the signal cluster arisen from \mathbf{I}_1 with $b_{i,1} = 1$ for $i = 1, \dots, 500$. Other 500 rows are purely noise, in *Oracle Group B*, which can be viewed as being centered around the n -dimensional origin. Therefore, a true proportion of null features is $\pi_0 = .50$.

We simulated three scenarios using $\sigma_b^2 = 5, 10, 15$ as an increasing noise level brings these two groups closer and makes the clustering task more difficult. PCA was applied on the dataset realized from each configuration to visualize the top 2 PCs (Figure S1). Being blind to *Oracle Groups*, the K -means clustering and the jackstraw tests were applied. Theoretically, the null p-values from the features that are not related to the latent variables (corresponding to *Oracle Group B*) should form the Uniform(0,1) distribution, which can be evaluated by the Kolmogorov-Smirnov (KS) test. We repeated a given simulation configuration 100 times independently and investigated how 100 KS test p-values from 100 independent simulations meet the joint null criterion [14].

We describe one simulation from the main scenario involving a moderate amount of noise $\sigma_b^2 = 10$. While 1000 features were split equally between *Cluster 1* and *2*, 30 and 470 null features were members of *Cluster 1* and *2*, respectively. Because *Cluster 1* contained 470 features related to the latent variable \mathbf{I}_1 , its center and \mathbf{I}_1 were highly correlated with a Pearson correlation of 0.99. The jackstraw test was then applied on the simulated data with $s = 100$ synthetic null features over $B = 5000$ iterations, while being blind to simulation parameters. Figure 2(a) shows histograms of p-values stratified by *Oracle Groups* as parametrized by dichotomous coefficients in \mathbf{B} . In *Oracle Group B*, the jackstraw p-values corresponding to 500 null features are uniformly distributed between 0 and 1. In contrast, the naive significance tests are highly anti-conservative, pushing towards 0. In *Oracle Group A*, the jackstraw p-values are greater than the naive p-values because the jackstraw approach learns the overfitting characteristics and fixes an anti-conservative bias (Figure 2(a)). Utilizing all m p-values, the proportion of null features are estimated to be $\widehat{\pi}_0 = 0.55$ for the jackstraw and $\widehat{\pi}_0 = 0.29$ for the naive methods.

We repeated this configuration to ensure accuracy and robustness across 100 independent simulations. In each simulation, we examined the joint behavior of 500 null p-values from *Oracle Group B* using a double-sided KS test. When the joint behavior of those KS test p-values follows the i.i.d. Uniform(0,1) distribution (where the double KS test p-value $> \alpha_{jnc}$), the subsequent multiple hypothesis testing procedures, including false discovery rates, hold true [14]. In other words, meeting the stringent standard of the joint null criterion demonstrates that the proposed methods overcome “double-dipping” inherent in utilizing cluster centers and membership assignments and that the p-values are jointly and marginally accurate [14]. A set of 100 KS test p-values, estimated from both the jackstraw and naive methods, are visualized against the Uniform(0,1) distribution (Figure 2(b)). The jackstraw tests satisfy the joint null criterion, where 100 KS test p-values are uniformly distributed (double KS test p-value = 0.79). In contrast, the naive methods are strongly anti-conservative, where 100 KS test p-values are strongly skewed towards 0 (double KS test

p-values $< 2.2 \times 10^{-16}$).

Results from two other simulation configurations, that are also independently repeated 100 times, are shown in Figure S2 and Figure S3. Simulated data with a relatively small noise $\sigma_b^2 = 5$ can be almost perfectly clustered. Nonetheless, the naive methods exhibit substantial overfitting where the double KS test p-value is $< 2.2 \times 10^{-16}$. The double KS test p-value for the jackstraw tests in this configuration is 0.81 (Figure S2). On the other hand, a greater noise with $\sigma_b^2 = 15$ represents a situation whose members of different clusters are substantially overlapping (Figure S1). The jackstraw tests indeed satisfy the joint null criterion with the double KS test p-value of 0.67 (Figure S3). Additional simulation studies further confirm that unlike the naive methods that overfit and produce an anti-conservative bias (downward deviations from the diagonal line), the proposed methods take account for uncertainty in clustering and result in valid p-values that enable rigorous error control.

Genomic Applications

Microarray Data from Yeast Cell Cycle Experiments

Cell cycle in *Saccharomyces cerevisiae* and other organisms are traditionally known to progress through discrete stages, such as M, G1, S, and so on [45]. With the advent of the microarray, gene expression levels from synchronized *S. cerevisiae* samples had been measured and analyzed in order to identify comprehensive sets of genes under cell cycle [8, 46–48]. However, in these conventional studies, experimentally verified genes under cell cycle are used to identify related genes that follow similar patterns. In contrast, we re-analyzed the expression data of 5981 genes from [8] in an unsupervised manner.

Genome-wide mRNA levels of elutriation-synchronized yeast cells were measured at 30 min intervals for 390 min (approximately 1 cell cycle) [8]. We processed and normalized this gene expression data according to [13, 49]. The number of clusters $K = 6$ was determined by the prior knowledge that there exist 6 stages of cell cycle [45]. While there are on-going debates on how to characterize and categorize cell cycle progression, $K = 6$ seems to be a reasonable choice. After having gotten $K = 6$ clusters from applying K -means clustering, we conducted the proposed jackstraw tests with $s = 300$ and $B = 10000$ to identify canonical genes within those clusters.

Histograms of p-values are shown in Figure 3(a), where the proportions of null features π_0 for 6 clusters are estimated to be .143, .149, .116, .178, .170, .175, .087, respectively. Note that the numeric values identifying those clusters are arbitrary without a meaningful order, but consistent within this manuscript. From a set of p-values, we calculated posterior probabilities that those genes are truly members of their assigned clusters. For example, among 709 genes that are originally assigned to the cluster 4, 45.1% (320) have posterior inclusion probabilities (PIPs) > 0.9 . Repeating this analysis for all 6 clusters, a total of 3826 genes are found to be significant at the same PIP threshold (Figure 3(c)). In other words, these are the canonical genes that drive the clusters of cell cycle.

Single Cell RNA-Seq Data from Jurkat and 293T Cells

Whereas conventional microarray and RNA-seq experiments obtain “bulk” gene expression from a sample that contains multiple cells, scRNA-seq enable more precise and accurate quantification from single cell samples. Recent studies using high-throughput and efficient scRNA-seq often measure gene expression from unlabeled single cells, in order to elucidate detailed molecular landscapes and identify cell identities (e.g., blood sub-types, sub-classifications of a disease) [10–12]. Commonly, the cell identities are determined by applying the clustering algorithms to their gene expression profiles.

We analyzed the scRNA-seq data from [12] that used a mixture of Jurkat and 293T cells (50:50). Note that while the mixture proportion is known, the identities of individual cells that have been sequenced are unknown. Because Jurkat (male and expressing *CD3D*) and 293T (female and expressing *XIST*) cell lines are highly distinct, we observed intelligible two groups separated along the 1st PC from their gene expression profiles [12]. However, massively parallel scRNA-seq regrettably generates multiplets (doublets, triplets, etc).

The rate of multiplets increases linearly with the recovered cell number, and through single nucleotide variant (SNV) detection, they inferred a 3.1% multiplet rate for this mixture experiment [12]. For ~ 10000 single cells, [12] reports $> 8\%$ multiplet rates. The ambiguous identities of single cells would become increasingly challenging as scRNA-seq becomes more affordable and widespread.

Following the original analysis pipeline, we applied the K -means clustering on the top 10 PCs based on unique molecular identifier (UMI) counts. The jackstraw tests for those $K = 2$ clusters were conducted with $s = 100$ and $B = 10000$. We found that the jackstraw p-values capture deviation away from two centers, along the 1st PC axis (Figure 4(a)). Using the q-value methodology [50], the proportion of null features (that are not members of the clusters) is estimated to be $\widehat{\pi}_0 = 0.05$. Then, we computed the proposed PIPs from p-values (Figure 4(b)). At $\text{PIP} < 0.80$ (equivalent to 20% local FDRs), 3.3% of 3381 single cells would be removed from corresponding clusters, effectively and automatically removing the majority of suspected multiplets. Instead of hard-thresholding the single cell samples at an arbitrary threshold, we can also visualize posterior probabilities as levels of transparency in a conventional PCA projection, where the top 2 PCs are plotted (Figure 4(c)). Please note that because dimension reduction does not fully capture local and global structures in the original high dimensions, distances in reduced dimensions (PCA, t-SNE, and alike) should be considered with caution.

Discussion

The explosion of biological data has increased the importance of unsupervised learning. Without the external and accurate labels for the observed data, unsupervised learning aims to estimate latent structure, reduce dimensions, and classify data features. In particular, clustering of high-dimensional genomic data has led to better understanding of and informative hypotheses for biological functions [8, 51], molecular subtypes [7, 52], and cell identities [53]. However, data-dependent classification cannot be used in downstream analyses without incurring spurious statistical significance. Our proposed methods solve this challenge by learning the uncertainty inherent in deriving clusters from the data and conducting a statistical test using the jackstraw strategy.

There exists a wide range of clustering algorithms to automatically assign the m observed features into K clusters. The proposed methods test whether the observed features are correctly assigned to the corresponding clusters. Our key ingredient is to generate and re-cluster the jackstraw data, which include a very small number s of synthetic null features. Because of $s \ll m$, the majority of observed features are intact, resulting in cluster centers that are almost identical to the original cluster centers. Subsequently, eventual assignments of s synthetic null features into K clusters are used to derive the empirical null distribution. We have demonstrated favorable operating characteristics using simulated and real genomic data. The proposed PIP methods open new possibilities for selecting canonical cluster members, shrinking cluster centers and improving cluster algorithms. Furthermore, the proposed methods may adaptively guide the choice of stable clusters.

Our proposed strategy enables rigorous application of unsupervised learning, such that the estimated latent structure can be re-used in downstream analyses. The jackstraw test for PCA and related methods [13] have been used in many specialized areas of genomic studies [11, 12, 51, 54–56]. Complementing this successful approach, we have developed the jackstraw test for clustering. It may be useful to integrate both variants of the jackstraw tests, from selecting highly informative genes to deriving cell identities. Differential expression analyses based on cluster-based cell identities may become more robust by incorporating the jackstraw tests. Because the proposed methods are not limited to genomics, we anticipate its adaptation in other fields of data-intensive science.

Acknowledgments

This research was supported in part by the Polish National Science Centre (NCN) grants 2016/23/D/ST6/03613.

Figures

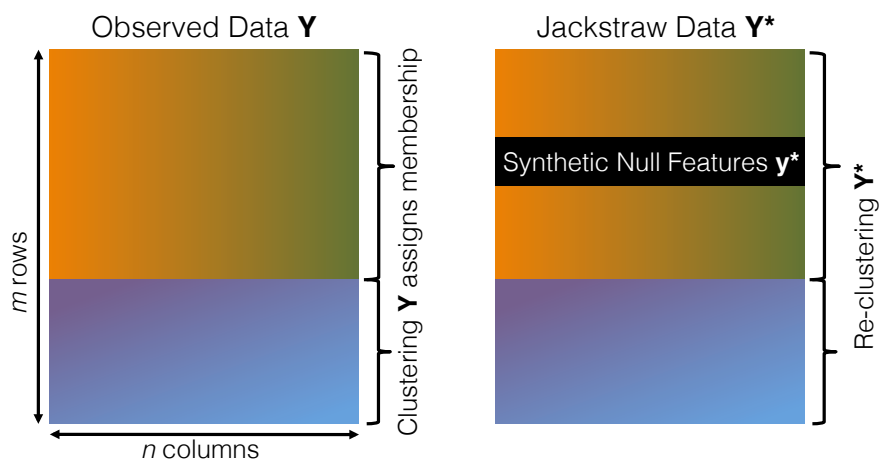


Figure 1: Illustration of the jackstraw data \mathbf{Y}^* consisted of s synthetic null features (\mathbf{y}^*) and $m - s$ observed features. By re-clustering \mathbf{Y}^* , \mathbf{y}^* that have been independently resampled are assigned into clusters. The jackstraw methods leverage this information to learn overfitting characteristics in clustering and to construct an empirical distribution of null statistics.

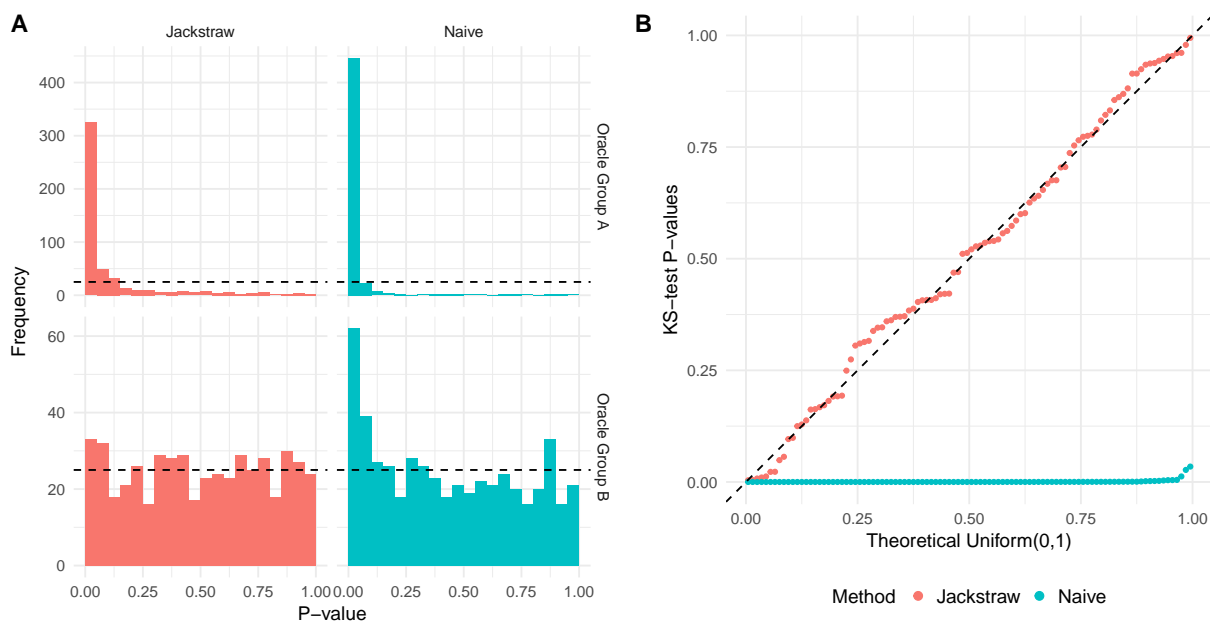


Figure 2: Evaluation of the jackstraw tests for clustering using the main simulation study with $\sigma = 10$. *Oracle Group A* contains 500 features that are derived from a latent variable $\mathbf{1} \stackrel{i.i.d}{\sim} \text{Normal}(0,1)$, whereas 500 features in *Oracle Group B* are noise. (a) Histograms of p-values stratified by methods. The jackstraw tests ($s = 100$, $B = 5000$) and the naive tests (at the same resolution) were applied without using any prior information. (b) This simulation study was repeated 100 times, where null p-values corresponding to *Oracle Group B* were evaluated against $\text{Uniform}(0,1)$ distribution using Kolmogorov-Smirnov (KS) tests. QQ plot of KS p-values from the jackstraw and naive methods are shown, where valid p-values follow a diagonal line. The double KS test p-values for the proposed jackstraw and naive methods are $= 0.79$ and $< 2.2 \times 10^{-16}$, respectively.

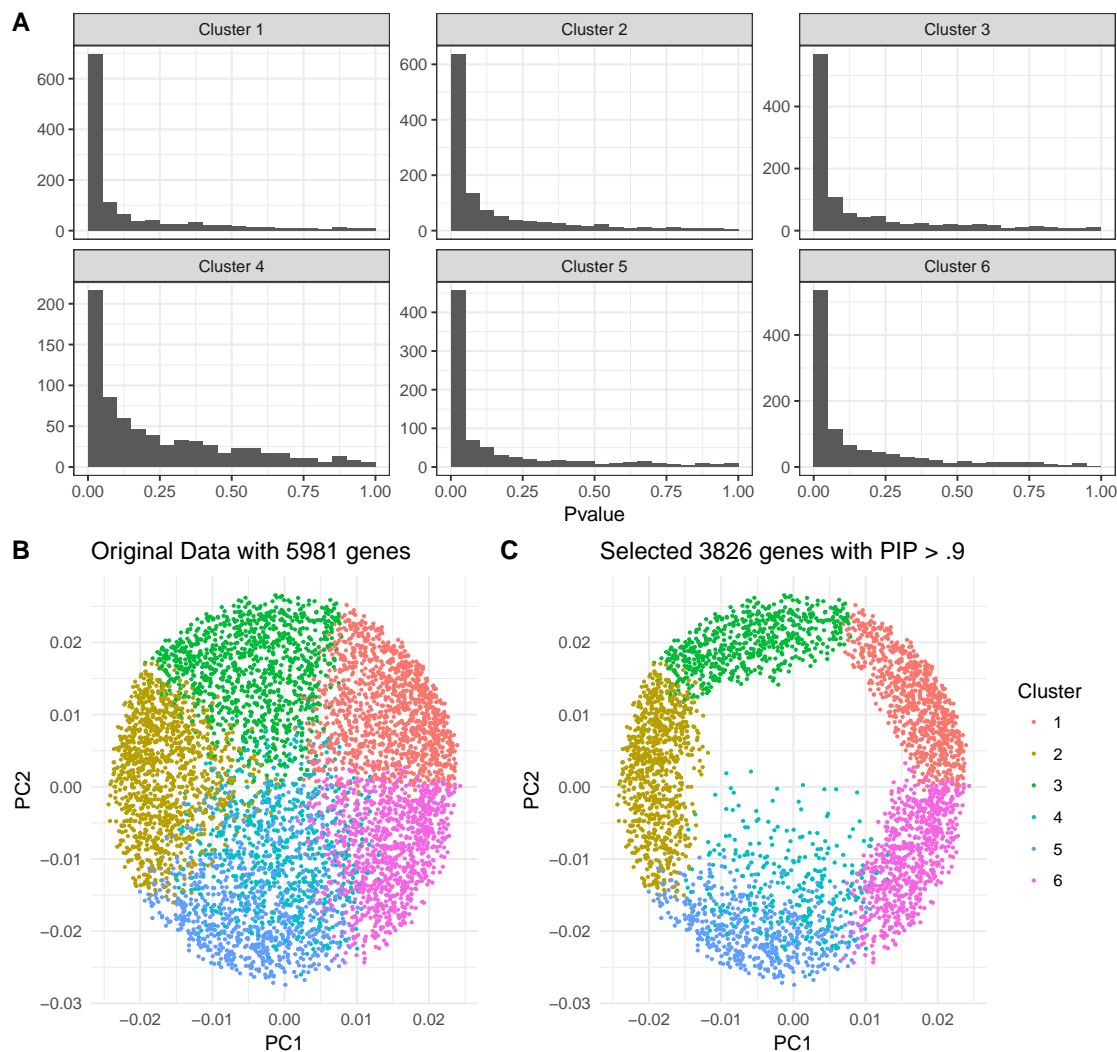


Figure 3: The jackstraw clustering analysis of microarray data of yeast cell cycle experiments. K -means clustering is applied on 5981 genes from [8] with $K = 6$. (a) Histograms of p-values from the yeast cell cycle gene expression profiles. The jackstraw tests for these 6 clusters are conducted with $s = 300$ and $B = 10000$. (b) The top 2 PC projection using the original data with 5981 genes. (c) The top 2 PC projection using 3826 genes with $PIP > .9$ from the proposed methods.

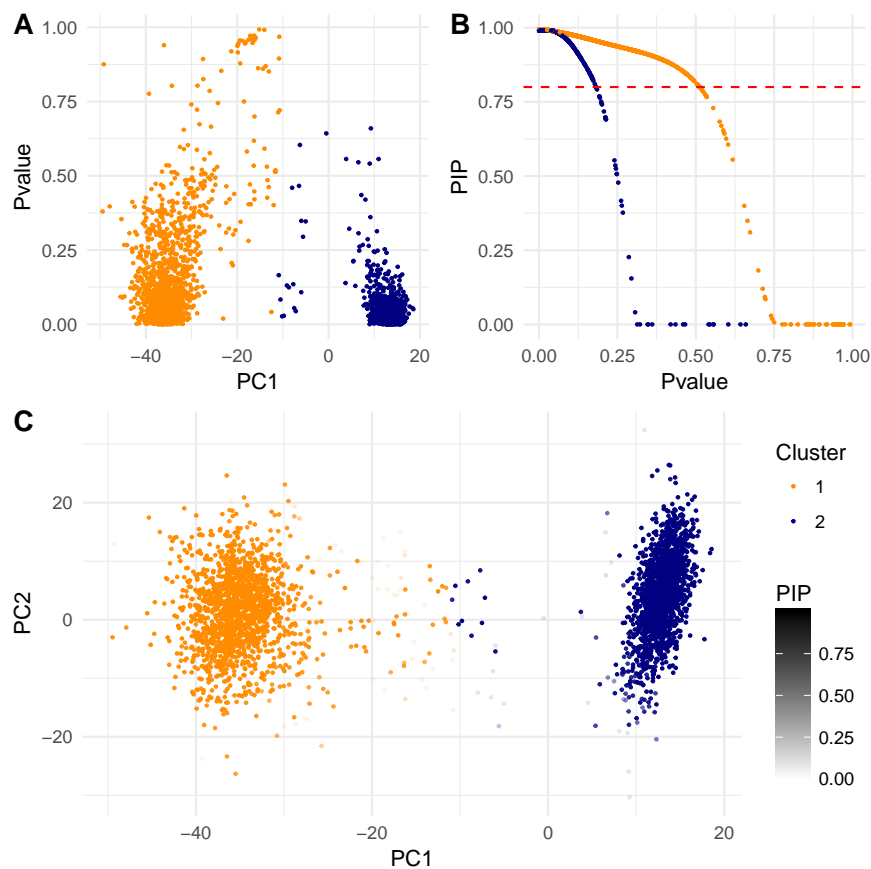


Figure 4: The jackstraw clustering analysis of scRNA-seq data of Jurkat and 293T cells. Two distinct cell lines are mixed at equal proportions (50:50) and sequenced without labeling using s10X Genomics' GemCode [12]. (a) Clustering membership p-values are plotted against the 1st PC, which largely separates two cell lines. (b) Posterior inclusion probabilities (PIPs) are computed from p-values. At a PIP threshold of 0.80, 3.3% of 3381 single cells would be discharged from corresponding clusters. (c) PIPs are visualized as alpha levels on the scatterplot of the top 2 PCs. Essentially, this is identical to Figure 2(e) in [12], except transparencies are set to PIPs from the proposed methods. Note that when PIP=0, as appeared in Figure 4(b), the data point is completely transparent.

Supplementary Figures

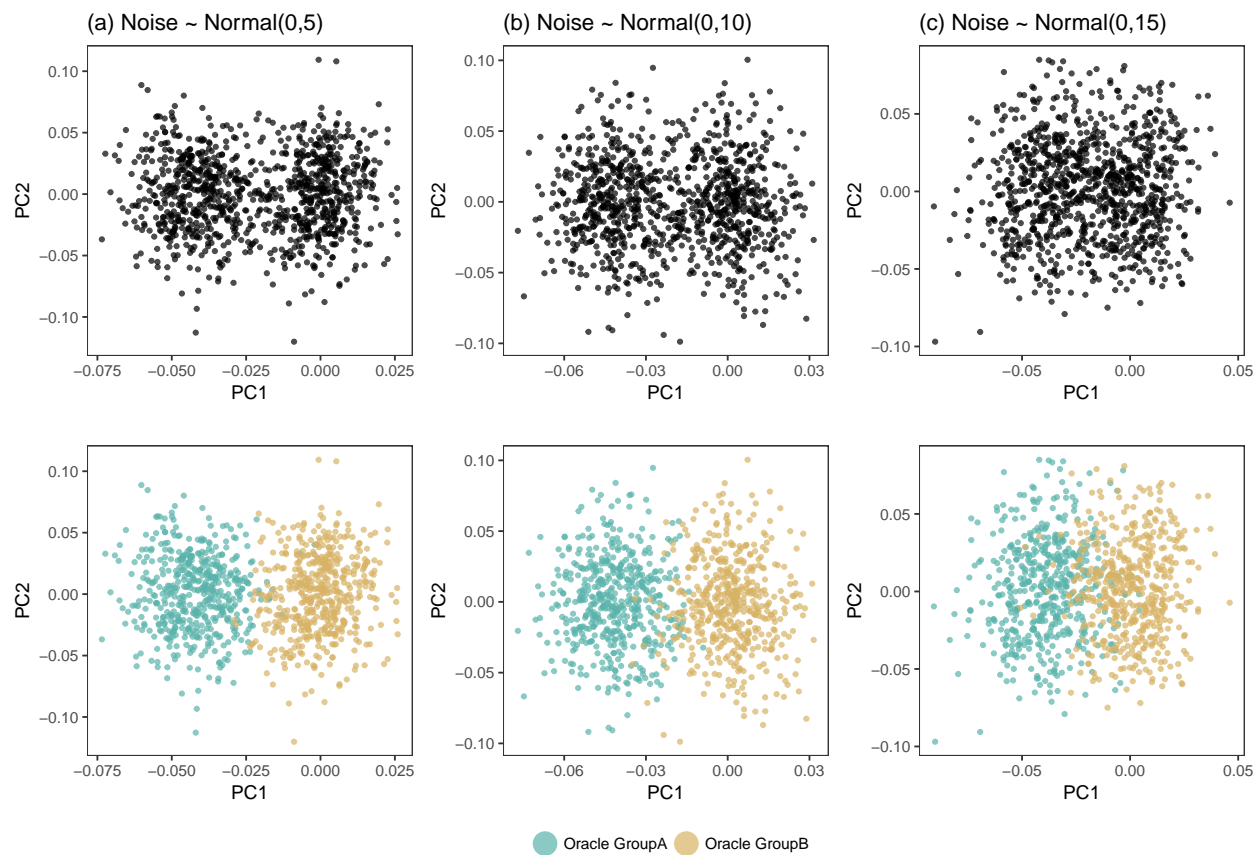


Figure S1: Scatterplots of the top 2 principal components (PCs) from the simulated data. *Oracle Groups* are shown in colors. An increasing level of noise, $\sigma^2 = 5, 10, 15$ brings data features from two different underlying structures closer together.

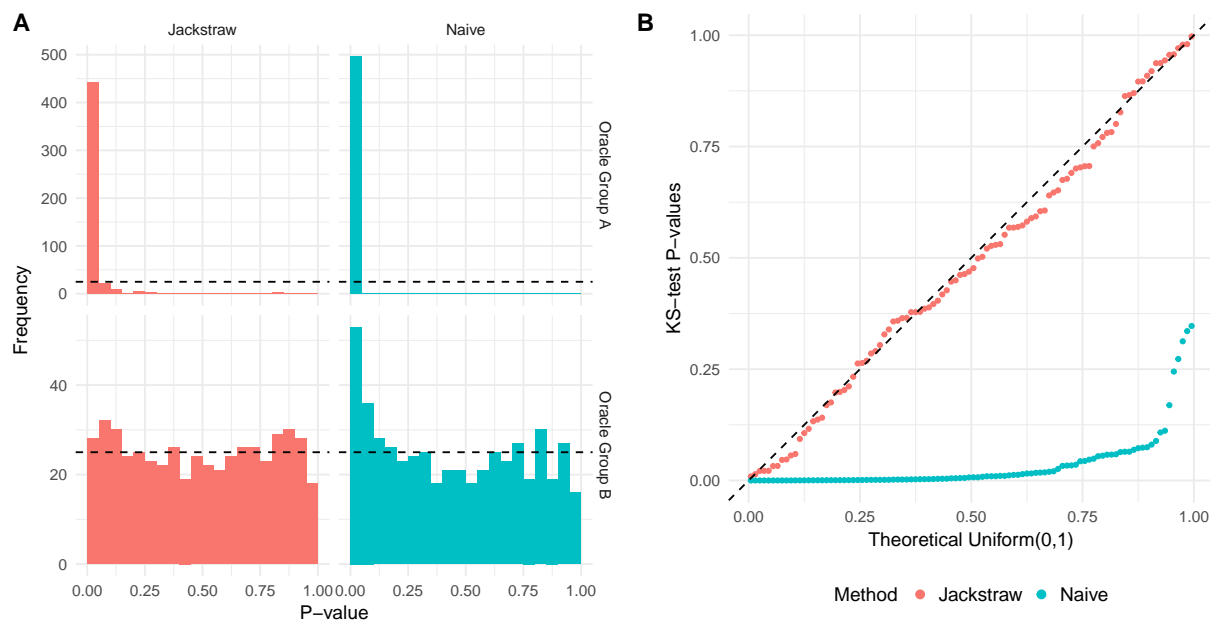


Figure S2: Simulation studies using $\sigma^2 = 5$. The jackstraw tests ($s = 100$ and $B = 5000$) or the naive tests are applied without using any information from simulation. (a) P-values are shown stratified by *Oracle Groups*, where the naive tests result in an anti-conservative bias. The uniformity of null p-values corresponding to *Oracle Group B* is examined by KS tests, which are independently repeated 100 times. (b) The total of 100 independent simulation studies are conducted, and 100 KS-test p-values are plotted against the Uniform(0,1) distribution. The proposed jackstraw tests meet the joint null criterion with a double KS test p-value of 0.81, whereas the naive tests are highly anti-conservative with a double KS test p-value of $< 2.2 \times 10^{-16}$.

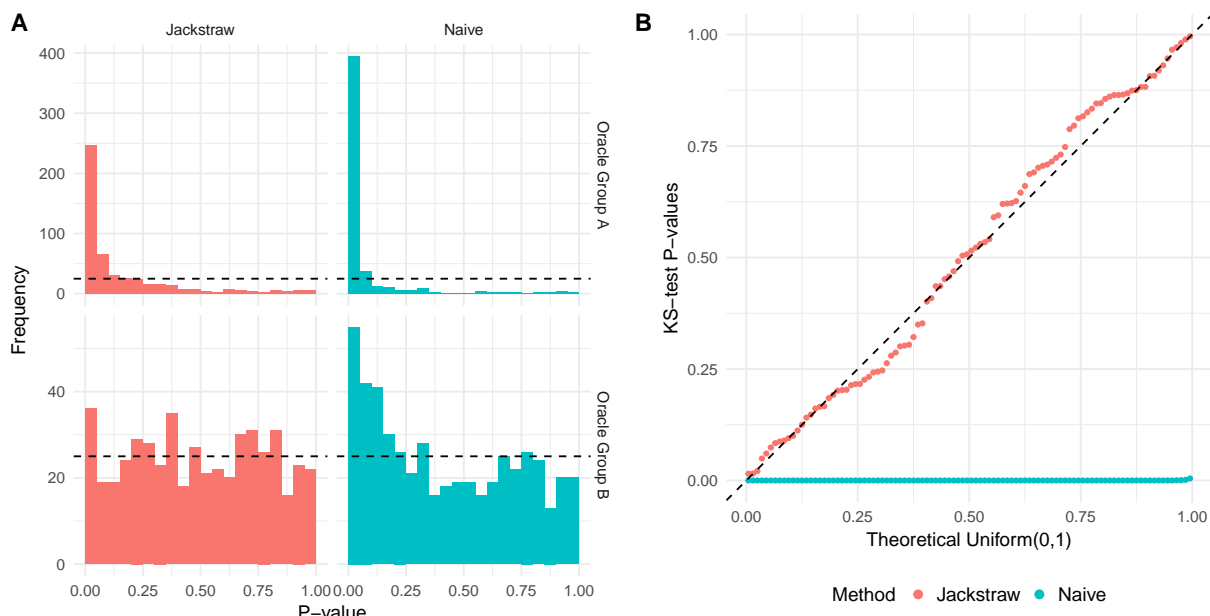


Figure S3: Simulation studies using $\sigma^2 = 15$. The jackstraw tests ($s = 100$ and $B = 5000$) or the naive tests are applied without using any information from simulation. (a) P-values are shown stratified by *Oracle Groups*, where the naive tests result in an anti-conservative bias. The uniformity of null p-values corresponding to *Oracle Group B* is examined by KS tests, which are independently repeated 100 times. (b) The total of 100 independent simulation studies are conducted, and 100 KS-test p-values are plotted against the Uniform(0,1) distribution. The proposed jackstraw tests meet the joint null criterion with a double KS test p-value of 0.67, whereas the naive tests are highly anti-conservative with a double KS test p-value of $< 2.2 \times 10^{-16}$.

References

1. Shalek, A. K., Satija, R., Adiconis, X., Gertner, R. S., Gaublomme, J. T., Raychowdhury, R., Schwartz, S., Yosef, N., Malboeuf, C., Lu, D., et al. Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells. *Nature* **498**(7453), 236–240 (2013).
2. Wills, Q. F., Livak, K. J., Tipping, A. J., Enver, T., Goldson, A. J., Sexton, D. W., and Holmes, C. Single-cell gene expression analysis reveals genetic associations masked in whole-tissue experiments. *Nature biotechnology* **31**(8), 748–752 (2013).
3. Rubakhin, S. S., Romanova, E. V., Nemes, P., and Sweedler, J. V. Profiling metabolites and peptides in single cells. *Nature methods* **8**(4s), S20–S29 (2011).
4. Budnik, B., Levy, E., and Slavov, N. Mass-spectrometry of single mammalian cells quantifies proteome heterogeneity during cell differentiation. *bioRxiv* **bioRxiv** (2017).
5. Alon, U., Barkai, N., Notterman, D. A., Gish, K., Ybarra, S., Mack, D., and Levine, A. J. Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proceedings of the National Academy of Sciences* **96**(12), 6745–6750 (1999).
6. Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J. P., Coller, H., Loh, M. L., Downing, J. R., Caligiuri, M. A., et al. Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *science* **286**(5439), 531–537 (1999).
7. Sørlie, T., Perou, C. M., Tibshirani, R., Aas, T., Geisler, S., Johnsen, H., Hastie, T., Eisen, M. B., Van De Rijn, M., Jeffrey, S. S., et al. Gene expression patterns of breast carcinomas distinguish tumor subclasses with clinical implications. *Proceedings of the National Academy of Sciences* **98**(19), 10869–10874 (2001).
8. Spellman, P. T., Sherlock, G., Zhang, M. Q., Iyer, V. R., Anders, K., Eisen, M. B., Brown, P. O., Botstein, D., and Futcher, B. Comprehensive identification of cell cycle-regulated genes of the yeast *saccharomyces cerevisiae* by microarray hybridization. *Molecular Biology of the Cell* **9**(12), 3273–3297 (1998).
9. Eisen, M. B., Spellman, P. T., Brown, P. O., and Botstein, D. Cluster analysis and display of genome-wide expression patterns. *Proceedings of the National Academy of Sciences of the United States of America* **95**(25), 14863–14868 (1998).
10. Jaitin, D. A., Kenigsberg, E., Keren-Shaul, H., Elefant, N., Paul, F., Zaretsky, I., Mildner, A., Cohen, N., Jung, S., Tanay, A., et al. Massively parallel single-cell rna-seq for marker-free decomposition of tissues into cell types. *Science* **343**(6172), 776–779 (2014).
11. Macosko, E. Z., Basu, A., Satija, R., Nemes, J., Shekhar, K., Goldman, M., Tirosh, I., Bialas, A. R., Kamitaki, N., Martersteck, E. M., et al. Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell* **161**(5), 1202–1214 (2015).
12. Zheng, G. X., Terry, J. M., Belgrader, P., Ryvkin, P., Bent, Z. W., Wilson, R., Ziraldo, S. B., Wheeler, T. D., McDermott, G. P., Zhu, J., et al. Massively parallel digital transcriptional profiling of single cells. *Nature communications* **8**, 14049 (2017).
13. Chung, N. C. and Storey, J. D. Statistical significance of variables driving systematic variation in high-dimensional data. *Bioinformatics* **31**(4), 545–554 (2015).
14. Leek, J. T. and Storey, J. D. The joint null criterion for multiple hypothesis tests. *Statistical Applications in Genetics and Molecular Biology* **10**(1), Article 28 (2011).
15. Linda M. Collins, S. T. L. *Latent Class and Latent Transition Analysis: With Applications in the Social, Behavioral, and Health Sciences*. John Wiley & Sons, (2010).

16. Bartholomew, D. J., Knott, M., and Moustaki, I. *Latent Variable Models and Factor Analysis: A Unified Approach*. Wiley Series in Probability and Statistics, (2011).
17. Yeung, K. Y., Fraley, C., Murua, A., Raftery, A. E., and Ruzzo, W. L. Model-based clustering and data transformations for gene expression data. *Bioinformatics* **17**(10), 977–987 (2001).
18. McLachlan, G. and Peel, D. *Finite mixture models*. John Wiley & Sons, (2004).
19. Fraley, C. and Raftery, A. E. Model-based methods of classification: using the mclust software in chemometrics. *Journal of Statistical Software* **18**(6), 1–13 (2007).
20. Thompson, E. A. and Geyer, C. J. Fuzzy p-values in latent variable problems. *Biometrika* **94**(1), 49–60 (2007).
21. Perkins, W., Tygert, M., and Ward, R. Significance testing without truth. *ArXiv stat.ME* (2013).
22. Meng, X.-L. Posterior predictive p-values. *The Annals of Statistics* **22**, 1142–1160 (1994).
23. Gelman, A., Meng, X.-L., and Stern, H. Posterior predictive assessment of model fitness via realized discrepancies. *Statistica sinica* **6**, 733–760 (1996).
24. Witten, D. M. and Tibshirani, R. A framework for feature selection in clustering. *Journal of the American Statistical Association* **105**(490), 713–726 (2010).
25. Sun, W., Wang, J., Fang, Y., et al. Regularized k-means clustering of high-dimensional data and its asymptotic consistency. *Electronic Journal of Statistics* **6**, 148–167 (2012).
26. Akaike, H. A new look at the statistical model identification. *IEEE transactions on automatic control* **19**(6), 716–723 (1974).
27. Schwarz, G. et al. Estimating the dimension of a model. *The annals of statistics* **6**(2), 461–464 (1978).
28. Bock, H.-H. On some significance tests in cluster analysis. *Journal of classification* **2**(1), 77–108 (1985).
29. Fraley, C. and Raftery, A. E. How many clusters? which clustering method? answers via model-based cluster analysis. *Computer Journal* **41**(8), 578–588 (1998).
30. Pelleg, D., Moore, A. W., et al. X-means: Extending k-means with efficient estimation of the number of clusters. In *ICML*, volume 1, 727–734, (2000).
31. Tibshirani, R., Walther, G., and Hastie, T. Estimating the number of clusters in a data set via the gap statistic. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **63**(2), 411–423 (2001).
32. Hamerly, G. and Elkan, C. Learning the k in k-means. In *Advances in neural information processing systems*, 281–288, (2004).
33. Liu, Y., Hayes, D. N., Nobel, A., and Marron, J. Statistical significance of clustering for high-dimension, low-sample size data. *Journal of the American Statistical Association* **103**(483), 1281–1293 (2008).
34. Huang, H., Liu, Y., Yuan, M., and Marron, J. Statistical significance of clustering using soft thresholding. *Journal of Computational and Graphical Statistics* **24**(4), 975–993 (2015).
35. MacQueen, J. et al. Some methods for classification and analysis of multivariate observations. In *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1 of 14, 281–297. Oakland, CA, USA., (1967).
36. Hartigan, J. A. and Wong, M. A. Algorithm as 136: A k-means clustering algorithm. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* **28**(1), 100–108 (1979).
37. Lloyd, S. Least squares quantization in pcm. *IEEE transactions on information theory* **28**(2), 129–137 (1982).

38. Efron, B. *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*, volume 1. Cambridge University Press, (2012).
39. Chung, N. C. *Statistical Inference of Variables Driving Systematic Variation in High-Dimensional Biological Data*. PhD thesis, Princeton University, (2014).
40. Efron, B., Tibshirani, R., Storey, J. D., and Tusher, V. Empirical bayes analysis of a microarray experiment. *Journal of the American Statistical Association* **96**(456), 1151–1160 (2001).
41. Efron, B. Size, power and false discovery rates. *The Annals of Statistics* **35**(4), 1351–1377, aug (2007).
42. Barbieri, M. and Berger, J. Optimal predictive model selection. *Annals of Statistics* **32**, 870–897 (2004).
43. Scott, J. and Berger, J. An exploration of aspects of bayesian multiple testing. *Journal of Statistical Planning and Inference* **136**, 2144–2162 (2005).
44. Ghosh, D., Chen, C., and Raghunathan, T. The false discovery rate: a variable selection perspective. *Journal of Statistical Planning and Inference* **136**, 2668–2684 (2006).
45. Alberts, B. *Molecular biology of the cell*. Garland science, (2017).
46. Cho, R. J., Campbell, M. J., Winzeler, E. A., Steinmetz, L., Conway, A., Wodicka, L., Wolfsberg, T. G., Gabrielian, A. E., Landsman, D., Lockhart, D. J., and Davis, R. W. A genome-wide transcriptional analysis of the mitotic cell cycle. *Mol Cell* **2**(1), 65–73, Jul (1998).
47. Tu, B. P., Kudlicki, A., Rowicka, M., and McKnight, S. L. Logic of the yeast metabolic cycle: temporal compartmentalization of cellular processes. *Science* **310**(5751), 1152–1158, Nov (2005).
48. Rowicka, M., Kudlicki, A., Tu, B. P., and Otwinowski, Z. High-resolution timing of cell cycle-regulated gene expression. *Proc Natl Acad Sci U S A* **104**(43), 16892–16897, Oct (2007).
49. Alter, O., Brown, P. O., and Botstein, D. Singular value decomposition for genome-wide expression data processing and modeling. *Proceedings of the National Academy of Sciences of the United States of America* **97**(18), 10101–10106 (2000).
50. Storey, J. D. and Tibshirani, R. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences* **100**(16), 9440–9445 (2003).
51. Jang, S., Choubey, S., Furchtgott, L., Zou, L.-N., Doyle, A., Menon, V., Loew, E. B., Krostag, A.-R., Martinez, R. A., Madisen, L., Levi, B. P., and Ramanathan, S. Dynamics of embryonic stem cell differentiation inferred from single-cell transcriptomics show a series of transitions through discrete cell states. *eLife* **6**, mar (2017).
52. Wang, D. and Gu, J. Integrative clustering methods of multi-omics data for molecule-based cancer classifications. *Quantitative Biology* **4**(1), 58–67, mar (2016).
53. Wagner, A., Regev, A., and Yosef, N. Revealing the vectors of cellular identity with single-cell genomics. *Nature Biotechnology* **34**(11), 1145–1160, nov (2016).
54. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F., and Regev, A. Spatial reconstruction of single-cell gene expression data. *Nature biotechnology* **33**(5), 495–502 (2015).
55. Chung, N., Szyda, J., Frąszczak, M., and Project, . B. G. Population structure analysis of bull genomes of european and western ancestry. *Scientific Reports* **7**(40688) (2017).
56. Farré, P., Jones, M. J., Meaney, M. J., Emberly, E., Turecki, G., and Kobor, M. S. Concordant and discordant DNA methylation signatures of aging in human blood and brain. *Epigenetics & Chromatin* **8**(1), may (2015).