

SCIENTIFIC REPORTS



OPEN

Detection of correlated hidden factors from single cell transcriptomes using Iteratively Adjusted-SVA (IA-SVA)

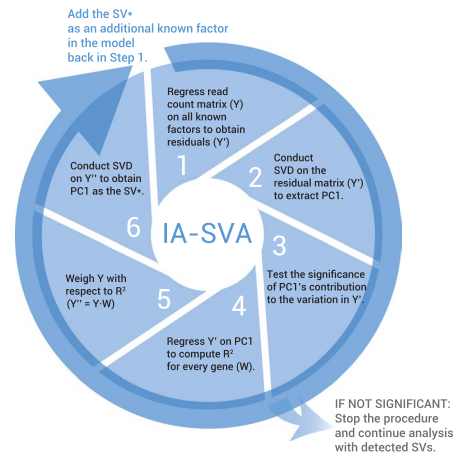
Donghyung Lee¹, Anthony Cheng^{1,2}, Nathan Lawlor¹, Mohan Bolisetty⁴ & Duygu Ucar^{1,2,3} 

Single cell RNA-sequencing (scRNA-seq) precisely characterizes gene expression levels and dissects variation in expression associated with the state (technical or biological) and the type of the cell, which is averaged out in bulk measurements. Multiple and correlated sources contribute to gene expression variation in single cells, which makes their estimation difficult with the existing methods developed for batch correction (e.g., surrogate variable analysis (SVA)) that estimate orthogonal transformations of these sources. We developed iteratively adjusted surrogate variable analysis (IA-SVA) that can estimate hidden factors even when they are correlated with other sources of variation by identifying a set of genes associated with each hidden factor in an iterative manner. Analysis of scRNA-seq data from human cells showed that IA-SVA could accurately capture hidden variation arising from technical (e.g., stacked doublet cells) or biological sources (e.g., cell type or cell-cycle stage). Furthermore, IA-SVA delivers a set of genes associated with the detected hidden source to be used in downstream data analyses. As a proof of concept, IA-SVA recapitulated known marker genes for islet cell subsets (e.g., alpha, beta), which improved the grouping of subsets into distinct clusters. Taken together, IA-SVA is an effective and novel method to dissect multiple and correlated sources of variation in scRNA-seq data.

Single-cell RNA-Sequencing (scRNA-seq) enables precise characterization of gene expression levels, which harbour variation in expression associated with both technical (e.g., biases in capturing transcripts from single cells, PCR amplifications or cell contamination) and biological sources (e.g., differences in cell cycle stage or cell types). If these sources are not accurately identified and properly accounted for, they might confound the downstream analyses and hence the biological conclusions^{1–3}. In bulk measurements, hidden sources of variation are typically ‘unwanted’ (e.g., batch effects) and are computationally eliminated from the data. However, in single cell RNA-seq data, variation/heterogeneity stemming from hidden biological sources can be the primary interest of the study; which necessitate their accurate detection (i.e., testing the existence of unknown heterogeneity in a cell population) and estimation (i.e., estimating a factor(s) representing the unknown heterogeneity (e.g., known cell subsets vs. unknown subset)) for downstream data analyses and interpretation. How hidden heterogeneity in single cell datasets can teach us novel biology was exemplified in a recent study that uncovered a rare subset of dendritic cells (DC), which only constitute 2–3% of the DC population⁴. Few genes were specifically expressed in this DC subset (e.g., AXL, SIGLEC1), which was captured by studying heterogeneity in single cell expression profiles that only affect a subset of genes and cells. This study exploited the variation in single cell expression profiles from blood samples to improve our knowledge of DC subsets. However, one challenge in detecting hidden sources of variation in scRNA-seq data lies in the existence of multiple and highly correlated hidden sources, including geometric library size (i.e., the total log-transformed read counts), number of expressed/detected genes in a cell, experimental batch effects, cell cycle stage and cell type^{5–8}. The correlated nature of hidden sources limits the efficacy of existing algorithms to accurately detect and estimate the source.

¹The Jackson Laboratory for Genomic Medicine, Farmington, 06032, CT, USA. ²Department of Genetics and Genome Sciences, University of Connecticut Health Center, Farmington, 06030, CT, USA. ³Institute of Systems Genomics, University of Connecticut Health Center, Farmington, 06030, CT, USA. ⁴Bristol-Myers Squibb, Pennington, NJ, 08534, USA. Correspondence and requests for materials should be addressed to D.L. (email: donghyung.lee@jax.org) or D.U. (email: duygu.ucar@jax.org)

A Iterative IA-SVA framework



B Application of IA-SVA on scRNA-seq data

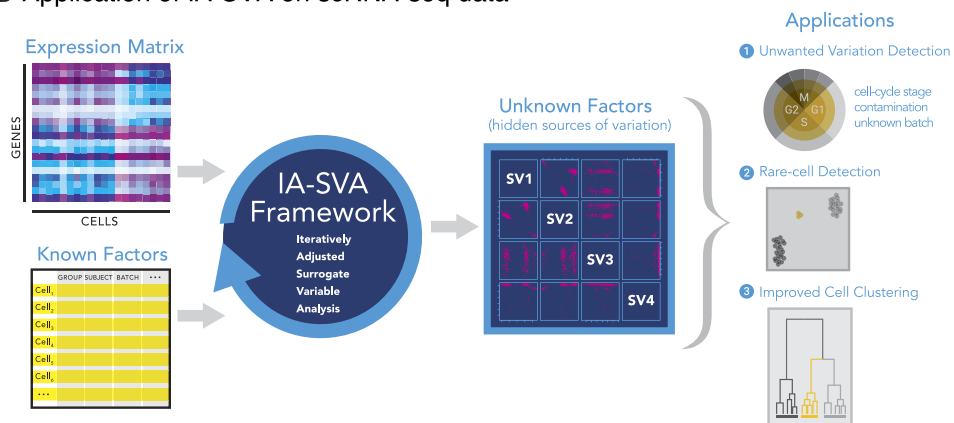


Figure 1. IA-SVA is a robust statistical framework to detect and estimate multiple and correlated hidden sources of variation. **(A)** Six-step IA-SVA procedure. IA-SVA computes the first principal component (PC1) from read counts adjusted for all known factors and tests its significance [Steps 1–3]. If significant, IA-SVA uses this PC1 to infer a set of genes associated with the hidden factor [Steps 4–5] and obtain a surrogate variable (SV) to represent the hidden factor using these genes [Step 6]. **(B)** IA-SVA uses single-cell gene expression data matrix and known factors to detect hidden sources of variation (e.g., cell contamination, cell-cycle status, and cell type). If these factors match to a biological variable of interest (e.g., cell type assignment), genes highly correlated with the factor can be detected and used in downstream analyses (e.g., data visualization).

‘Surrogate variable analysis’ (SVA)^{9–11} is a family of algorithms that are developed to detect and remove hidden “unwanted” variation (e.g., batch effect) in gene expression data by accurately parsing the data into signal and noise. A number of SVA-based methods have been developed and used for the analyses of microarray, bulk, and single-cell RNA-seq data including SSSVA¹¹ (supervised surrogate variable analysis), USVA¹⁰ (unsupervised SVA), ISVA¹² (Independent SVA), RUV (removing unwanted variation)^{13,14}, and most recently scLVM⁶ (single-cell latent variable model). These methods primarily aim to remove ‘unwanted’ variation (e.g., batch or cell-cycle effect) in data while preserving the biological signal of interest typically to improve downstream differential expression analyses between cases and controls. For this purpose, they utilize PCA (principal component analysis), SVD (singular value decomposition) or ICA (independent component analysis) to infer orthogonal transformations of hidden factors that can be used as covariates in downstream analysis. This paradigm by definition results in orthogonality between multiple estimated (and known) factors, which is a desired feature of batch correction methods in order to protect the signal of interest in downstream differential analysis¹⁴. However, this orthogonality assumption limits the efficacy of existing SVA-based methods to precisely estimate the sources of ‘wanted’ variation in single cell RNA-seq data.

To fill this gap, we developed Iteratively Adjusted Surrogate Variable Analysis (IA-SVA) (Fig. 1A and Methods for details) to estimate hidden factors even if these factors are correlated with known factors. IA-SVA provides three major advantages in single cell data analyses. First, it accurately estimates multiple hidden sources of variation even if the sources are correlated with each other and with known sources, which is not possible with existing SVA-based methods. Second, it enables assessing the significance of each detected factor for explaining the unmodeled variation in the data. Third, it identifies a set of genes that are significantly associated with the

	USVA	SSVA	IA-SVA	USVA	SSVA	IA-SVA
	r = 0.3 ~ 0.6			r < 0.3		
Power* (F1**)	1	1	1	1	1	1
Power (F2)	1	1	1	1	1	1
Power (F3)	0.78	0.78	0.87	1	1	1
Cor*** (F1)	0.93	0.95	0.95	0.98	0.98	1
Cor (F2)	0.72	0.75	0.94	0.94	0.94	0.99
Cor (F3)	0.75	0.78	0.95	0.93	0.93	0.98
	USVA	SSVA	IA-SVA			
Type I error*	0.09	0.09	0.04			

Table 1. IA-SVA accurately captures unknown sources of variation while controlling Type I error rate at a nominal level. Empirical power, Type I error rate, and the accuracy of estimates for IA-SVA, SSVA, and USVA assessed using simulated single-cell gene expression data. Alternative scenarios are simulated in which hidden factors are moderately ($|r| = \sim 0.3\text{--}0.6$, first three columns) or weakly ($|r| < 0.3$, last three columns) correlated with the group variable. IA-SVA outperforms alternative methods especially while detecting variation stemming from a smaller fraction of genes (10%) and especially when factors are correlated. *Nominal Type I error rate: 0.05. **F1, F2, F3 refers to Factor1, Factor2, and Factor3. ***Average of the absolute Pearson correlation coefficients between the true factor and the estimated factor is used as the accuracy measure.

detected hidden source. IA-SVA can be used to detect variation stemming from different sources in scRNA-seq data including cell contamination, cell cycle differences, or differences in cell types (Fig. 1B). In simulation studies we showed that IA-SVA (i) provides high statistical power in detecting hidden factors; (ii) controls Type I error rate at the nominal level ($\alpha = 0.05$); (iii) delivers high accuracy in estimating hidden factors. We evaluated the efficacy of IA-SVA on scRNA-seq data from human pancreatic islets and brain cells and showed that IA-SVA is effective in capturing heterogeneity associated with both technical (e.g., doublet cells) and biological sources (e.g., differences in cell types or cell-cycle stages). Furthermore, we showed that IA-SVA based gene selection can be further utilized in downstream analyses such as data visualization using t-distributed stochastic neighbor embedding (tSNE)¹⁵ and performs favourably compared to existing methods developed for gene selection and visualization (e.g., Spectral tSNE¹⁶).

Results

Benchmarking IA-SVA on simulated and real data. To assess and compare the detection power, Type I error rate, and the accuracy of hidden source estimates using IA-SVA and existing state-of-the-art methods (i.e., USVA and SSVA), we performed simulation studies (see Methods for details) under the null hypothesis (i.e., a group (e.g., males vs. females) variable affecting 10% of genes and no hidden factor) and under the alternative hypothesis (i.e., a group variable and three hidden factors affecting 10%, 30%, 20% and 10% of genes, respectively). Under the alternative hypothesis, we considered two correlation scenarios where the three hidden factors are moderately ($|r| = \sim 0.3\text{--}0.6$) or weakly ($|r| < 0.3$) correlated with the group variable (i.e., a known factor). Under each simulation scenario, we generated 1,000 scRNA-seq data sets (10,000 genes and 50 cells) and applied IA-SVA, USVA and SSVA ($\alpha = 0.05$, 50 permutations) on them to detect simulated hidden factors. Using these simulation results, we assessed the empirical Type I error rate of each method (i.e., the number of times each method detects a false positive factor under the null hypothesis at the nominal level of 0.05 divided by the number of simulations ($n = 1,000$)). Similarly, we also quantified the empirical detection power rate of each method under different alternative hypothesis scenarios as the number of times each method detects a simulated factor under the alternative hypothesis (i.e., a factor actually exists and is detected as significant by the method) divided by the number of simulations. We used the average of the absolute Pearson correlation coefficients between the simulated and estimated hidden factors to quantify the accuracy of estimates.

Simulation studies showed that IA-SVA performs equally well or better than USVA and SSVA in terms of the detection power and the accuracy of the estimate while controlling the Type I error rate (0.04 for IA-SVA versus 0.09 for USVA and SSVA) (Table 1). In particular, IA-SVA was more effective when a hidden factor affected a small percentage of genes and when the factors were correlated ($|r| = 0.3\text{--}0.6$) with the known factor (i.e., group variable). For example, IA-SVA detected Factor3, which affected only 10% genes, 87% of the time, whereas USVA and SSVA detected this factor 78% of the simulations (first three columns in Table 1). More importantly, IA-SVA correctly inferred the correlations among multiple hidden factors while USVA and SSVA delivered biased estimates due to their orthogonality assumption (Supplementary Fig. S1).

We further compared IA-SVA against other existing algorithms including ZINB-WaVe¹⁷ - a method that is developed for single-cell RNA-seq data analyses and uses a zero inflated negative binomial model. For benchmarking purposes, we used scRNA-seq data from human cortex¹⁸ and used IA-SVA, PCA, SVA¹¹, RUV-g and RUV-r¹³, and ZINB-WaVe to uncover the surrogate variable (SV) associated with the cell type assignment (neuron vs. astrocyte). To make the task more challenging, we only utilized 1,000 randomly selected genes that are adequately expressed (read count >3 for at least 3 cells) and repeated the analyses 100 times to eliminate the sampling bias. For each method (except for PCA), patient identity and geometric library size were used as covariates. Among the top 3 SVs identified for each method, the SV that correlates the best with the known cell type assignments was used for evaluation purposes to eliminate any bias especially for PCA. These analyses showed that IA-SVA outperforms other algorithms in terms of its ability to accurately identify the SV that is associated

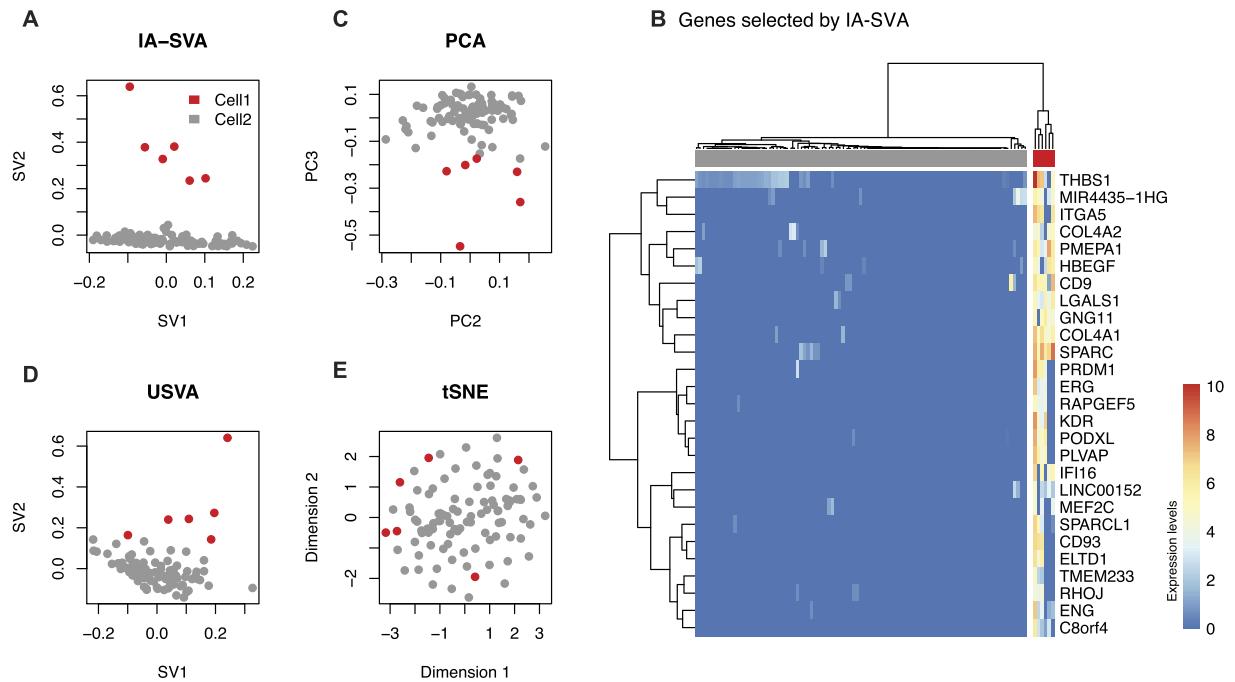


Figure 2. IA-SVA can detect heterogeneity originating from potentially contaminated alpha cells. (A) Outlier alpha cells captured using IA-SVA and same cells marked in respective (C) PCA, (D) USVA, and (E) tSNE analyses. Cells are clustered into two groups (red vs. gray dots) based on IA-SVA's surrogate variable 2 (SV2 > 0.1). In PCA, PC1 was discarded since it explains the geometric library size. (B) Hierarchical clustering of alpha cells using 27 genes significantly associated with SV2 (FDR < 0.05 and $R^2 > 0.6$) (ward.D2 and cutree_cols = 2). 6 cells clearly separated from the rest of the alpha cells based on the expression of these 27 genes. The values in the heatmap and the color bar are log-transformed (base e) normalized counts, i.e., $\log_e(1 + \text{read counts normalized using SCnorm})$.

with cell type assignment in human brain cells (see Supplementary Fig. S2A for a single run and S2B for the average of 100 runs). The SV inferred by IA-SVA was highly correlated with the true cell assignments compared to other methods (average correlation is 0.926 for 100 runs). SVA and PCA were the second-best performers for this task (average correlations around 0.85). ZINB-WaVe ($\text{cor} = 0.81$) was more effective than RUV-r, however, was the slowest among all tested algorithms (average run time 162 seconds vs. 12 seconds for IA-SVA using the permutation procedure with 20 repeats) (Supplementary Fig. S2). These results suggest that among the six tested methods, IA-SVA is the most effective in inferring the SV associated with the cell type assignments in single human brain cells.

IA-SVA captures variation stemming from a small number of alpha cells. To test whether IA-SVA is effective in capturing hidden variation within a homogenous cell population, we analysed scRNA-seq data generated from human alpha cells ($n = 101$, marked with glucagon (GCG) expression) obtained from three diabetic patients¹⁹ using the Fluidigm C1 platform²⁰, for which the original study did not report any separation of these alpha cells. Using geometric library size and patient ID as known factors, significant SVs were inferred using IA-SVA ($\alpha = 0.05$, 50 permutations) on the data (14,416 genes and 101 cells). For comparison, we applied PCA, USVA, and tSNE on this data. In USVA analysis, similarly geometric library size and patient ID were used as known factors and significant SVs were obtained ($\alpha = 0.05$, 50 permutations). In the PCA analysis, PC1 was discarded since it is highly correlated ($r = 0.99$) with the geometric library size.

The top two significant SVs inferred by IA-SVA clearly separated alpha cells into two groups (six outlier cells marked in red vs. the rest marked in grey at SV2 > 0.1) (Fig. 2A). 27 genes significantly associated with second SV (SV2) (Benjamini-Hochberg q-value (FDR) < 0.05, coefficient of determination (R^2) > 0.6), which included genes expressed in fibroblasts such as *COL4A1* and *COL4A2*. These genes were exclusively expressed in six outlier cells and clearly separated alpha cells into two clusters (Fig. 2B). A larger set of SV2-associated genes ($n = 108$, FDR < 0.05, $R^2 > 0.3$) was used for pathway and GO enrichment analyses and uncovered that these genes are associated with extracellular matrix receptors (Supplementary Table S1). Hence, these outlier cells likely arise from cell contamination (e.g., fibroblasts contaminating islet cells) or cell doublets (e.g., two cells captured together) — a known problem in early Fluidigm C1 experiments^{21,22}. Alternative methods (i.e., PCA, USVA, tSNE) failed to clearly detect these outlier cells (Fig. 2C–E).

We next studied whether this source of heterogeneity can be recapitulated in an independent and bigger human islet scRNA-seq dataset²⁰, using gene expression profiles (17,168 genes) of 569 alpha cells from six diabetic patients. Using geometric library size and patient ID as known factors we identified top 2 significant SVs using IA-SVA and USVA. For comparison, we also conducted PCA and tSNE analyses on this data. In PCA, PC1

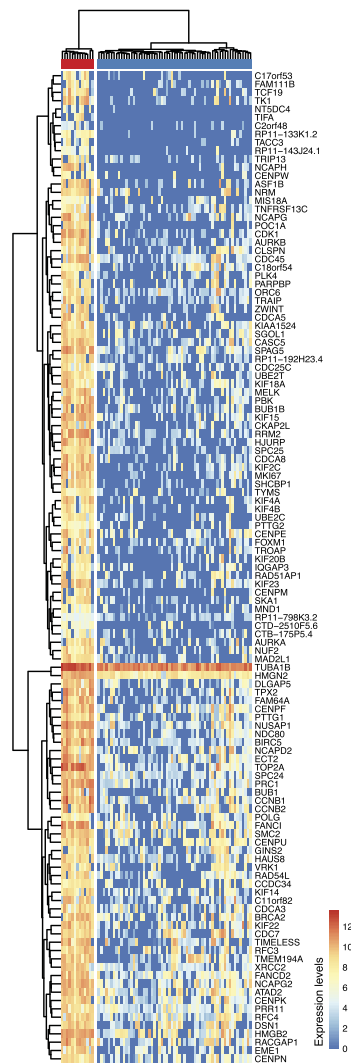
was discarded since it matched the geometric library size, which is adjusted for in IA-SVA and USVA analyses. IA-SVA's SV2 separated alpha cells into two groups (Supplementary Fig. S3A) and as in the previous case it was associated with fibrotic response genes including *SPARC*, *COL4A1*, *COL4A2* ($n = 81$, $FDR < 0.05$ and $R^2 > 0.3$) (Supplementary Fig. S3B, GO/pathway results in Supplementary Table S2). These results highlight IA-SVA's ability to detect variation among alpha cells potentially due to cell contamination or cell doublets. USVA and tSNE failed to clearly separate these compromised alpha cells (Supplementary Fig. S3D,E) from the rest of the cells. PCA's performance in separating these cells was more comparable to that of IA-SVA (Supplementary Fig. S3C), however, the clustering of cells was still confounded by known technical factors (i.e., Patient identity).

IA-SVA accurately detects variation arising from cell-cycle stage differences. Differences in cell-cycle stages lead to variation in single cell gene expression data³. Supervised methods based on SVA have been developed to detect and correct for cell cycle stage differences, most notably the scLVM algorithm. scLVM implements a Bayesian latent variable model to infer hidden cell-cycle factors by using known cell cycle genes⁶. IA-SVA can provide an unsupervised alternative by accurately capturing cell-cycle related variation in single cell data. To show this, we analyzed scRNA-seq data (21,907 genes and 74 cells) obtained from human glioblastomas that has an established cell-cycle signature²³. We conducted IA-SVA analyses by using geometric library size as a known factor and extracted top 2 significant SVs ($\alpha = 0.05$, 50 permutations).

IA-SVA's SV1 clearly separated 13 cells from the rest (cells marked in red in Fig. 3A), which was associated with 119 genes ($FDR < 0.05$ and $R^2 > 0.3$). Pathway and GO enrichment analyses of these genes^{24,25} revealed significant enrichment for cell-cycle process related GO terms and pathways (Fig. 3B, Supplementary Table S3), suggesting that this hidden variation is stemming from cell-cycle stage differences. Indeed, cell-cycle-stage predictions of cells using the SCRAN R package²⁶ showed that cells in different cell-cycle stages have different SV1 values (Fig. 3C). We noted that SV1 is highly correlated ($|r| = 0.44$) with the geometric library size (typically the top contributor to the variation in single cell data). These results demonstrate that IA-SVA can effectively detect variation stemming from cell-cycle differences in an unsupervised manner from single cell transcriptomes, even if this factor is highly correlated with known factors.

IA-SVA based gene selection improves single cell data visualization. tSNE and other dimension reduction algorithms (e.g., Spectral tSNE implemented in Seurat¹⁶) are frequently used to visualize single cell data since they group together cells with similar gene expression patterns. However, variation introduced by technical or biological factors can confound the signal of interest and generate spurious clustering of data. IA-SVA can be particularly effective in handling this problem by estimating hidden factors of interest accurately while adjusting for all known factors of no interest. Moreover, IA-SVA identifies genes associated with each detected hidden factor, which could be biologically relevant such as marker genes for different cell types. The genes inferred by IA-SVA can significantly improve the performance of data visualization methods (e.g., tSNE¹⁵). To illustrate this, we studied single cell gene expression profiles (16,005 genes) of alpha ($n = 101$, marked with glucagon (*GCG*) expression), beta ($n = 96$, marked with insulin (*INS*) expression), and ductal ($n = 16$, marked with *KRT19* expression) cells obtained from three diabetic patients¹⁹. First, we applied tSNE on all genes ($n = 16,005$) and color-coded genes based on the reported cell type assignments¹⁹, which failed to separate cells from different origins (Fig. 4A). Next, we applied IA-SVA on this data using patient ID, batch ID and the number of expressed genes as known factors and obtained significant SVs. SV1 and SV2 separated cells into distinct clusters (Supplementary Fig. S4), suggesting that these SVs might be associated with cell type differences. Indeed, genes associated with SV1 and SV2 ($n = 92$, $FDR < 0.05$ and $R^2 > 0.5$) included known marker genes used in the original study (*INS*, *GCG*, *KRT19*) and uncovered alternative marker genes associated with alpha, beta and ductal cells (Fig. 4B). These genes were annotated with diabetes and insulin processing related GO terms and pathways (Supplementary Table S4). As expected, tSNE analyses based on these 92 genes improved data visualization and clearly grouped together cells with respect to their cell type assignments (Fig. 4C). Such improved analyses can be instrumental in discovering cells that might be incorrectly labelled based on a single marker gene. For example, our analyses revealed a beta cell that is labelled as a ductal cell in the original study (one green cell clustered with blue cells in Fig. 4C). For comparison, we applied recently developed visualization methods, CellView²⁷ and Spectral tSNE¹⁶, on the same data with their recommended settings. CellView identified the 1000 most over-dispersed genes and conducted tSNE on these genes. Spectral tSNE detected 2,933 most over-dispersed genes and performed tSNE on significant principal components of these genes. On this small dataset, both methods managed to group cells of different types into distinct groups (Supplementary Fig. S5), suggesting that existing methods for gene selection and visualization are effective when datasets are small in size and are not confounded with multiple factors.

To test the efficacy of these methods on a bigger and more complex dataset, we conducted similar analyses on scRNA-seq data (19,226 genes) of 1,600 islet cells including alpha ($n = 946$), beta ($n = 503$), delta ($n = 58$), and PP ($n = 93$) cells from 6 diabetic and 12 non-diabetic individuals, where the study includes multiple confounding factors (e.g., ethnicity, disease state)²⁰. We noted that original cell type assignments significantly correlate with patient identifications ($C = 0.48$, $C =$ Pearson's contingency coefficient) and with ethnicity ($C = 0.25$), which would reduce the ability of existing methods to detect variation associated with cell types. In such complex datasets, failing to properly adjust for potential confounding factors prior to data analyses can lead to spurious grouping of cells, which might mislead the biological conclusions. Indeed, when these cells were visualized using tSNE using all genes ($n = 19,226$) and were color-coded with respect to the original cell-type assignments²⁰, cell types did not separate from each other and spurious clusters were observed within each cell type (Fig. 4D). As suspected, potential confounding factors (i.e., patient ID and ethnicity) explained this grouping of cells (Supplementary Fig. S6), which might be misleading as researchers are looking for alpha and beta cell subtypes that can be related to Type 2 Diabetes pathogenesis²⁸. To eliminate spurious clusters stemming from known

A Genes selected by IA-SVA**B** GO annotations for selected genes

GO ID	Term	logP	#Genes in term	#Genes in list	%Target
GO:0022402	cell cycle process	-171.95	1073	80	0.78
GO:0000280	nuclear division	-163.02	509	64	0.62
GO:0048285	organelle fission	-159.31	538	64	0.62
GO:0007049	cell cycle	-159.13	1317	81	0.79
GO:0000278	mitotic cell cycle	-153.29	781	69	0.67
GO:1903047	mitotic cell cycle process	-152.04	753	68	0.66
GO:0007067	mitotic nuclear division	-140.84	379	54	0.52
GO:0007059	chromosome segregation	-123.15	289	46	0.45
GO:0051301	cell division	-121.70	466	52	0.50
GO:0051276	chromosome organization	-121.55	570	55	0.53
GO:0098813	nuclear chromosome segreg.	-118.92	245	43	0.42
GO:0000819	sister chromatid segregation	-112.24	178	38	0.37
GO:0006996	organelle organization	-87.63	2799	77	0.75
GO:0000070	mitotic sister chromatid segreg.	-77.58	98	25	0.24
GO:0007062	sister chromatid cohesion	-65.51	113	23	0.22
GO:1902589	single-organism organelle organ.	-58.90	1380	49	0.47
GO:0016043	cellular component organization	-57.15	5162	83	0.80
GO:0071840	cellular comp. org. or biogenesis	-54.85	5323	83	0.80
GO:0010564	regulation of cell cycle process	-50.61	559	32	0.31

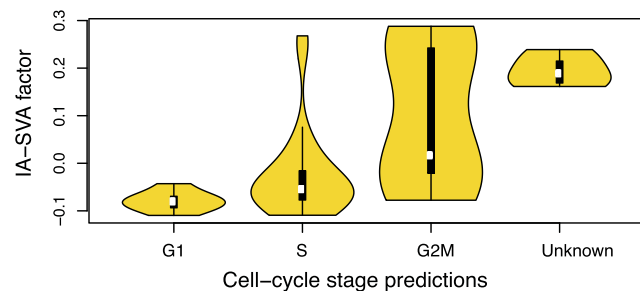
C IA-SVA factor vs. cell-cycle-stage predictions

Figure 3. IA-SVA can detect heterogeneity stemming from differences in cell-cycle stage. **(A)** Hierarchical clustering using 119 genes that are significantly associated ($FDR < 0.05$ and $R^2 > 0.3$) with IA-SVA's SV1. Clustering analysis confirms the separation of cells into two groups based on these genes (ward.D2 and `cutree_cols = 2`). The values in the heatmap and the color bar are log-transformed (base e) normalized counts, i.e., $\log_e(1 + \text{read counts normalized using SCnorm})$. **(B)** Top GO annotations for 119 selected genes. Majority of these genes were associated with cell-cycle related GO terms. **(C)** IA-SVA's SV1 can segregate cells based on their cell-cycle-stage as predicted by SCRAN.

factors, existing methods (e.g., Seurat) simply regress out all known factors prior to visualization. However, this might affect the signal of interest (i.e., cell type assignment), due to the high correlation between known factors (i.e., patient ID) and the hidden factor (i.e., cell types).

We applied IA-SVA on this complex data, while accounting for known factors (i.e., the number of expressed genes and patient ID) and extracted top four significant SVs (Supplementary Fig. S7A,B). We identified 57 genes associated with the most significant SV (SV1) ($FDR < 0.05$ and $R^2 > 0.5$), which included known marker genes (i.e., *INS* and *GCG*) (Supplementary Fig. S8, Table S5) and revealed novel marker genes for these cell types. tSNE analyses using these 57 IA-SVA detected genes clearly separated different cell types into discrete groups and reinforced the importance of properly adjusting for known factors prior to data analyses (Fig. 4E). For comparison, we applied CellView and Spectral tSNE on this data with recommended settings; however they failed to accurately group cells into distinct cell types (Fig. 4F,G). Similar analyses were conducted using PCA and USVA on the same data, where top surrogate factors obtained with both methods failed to separate different cell types into distinct groups (Supplementary Fig. S7C,D). Combined together these analyses suggest that IA-SVA is particularly effective in the analyses of complex datasets, which include the measurements of many cells that are affected by diverse confounding factors.

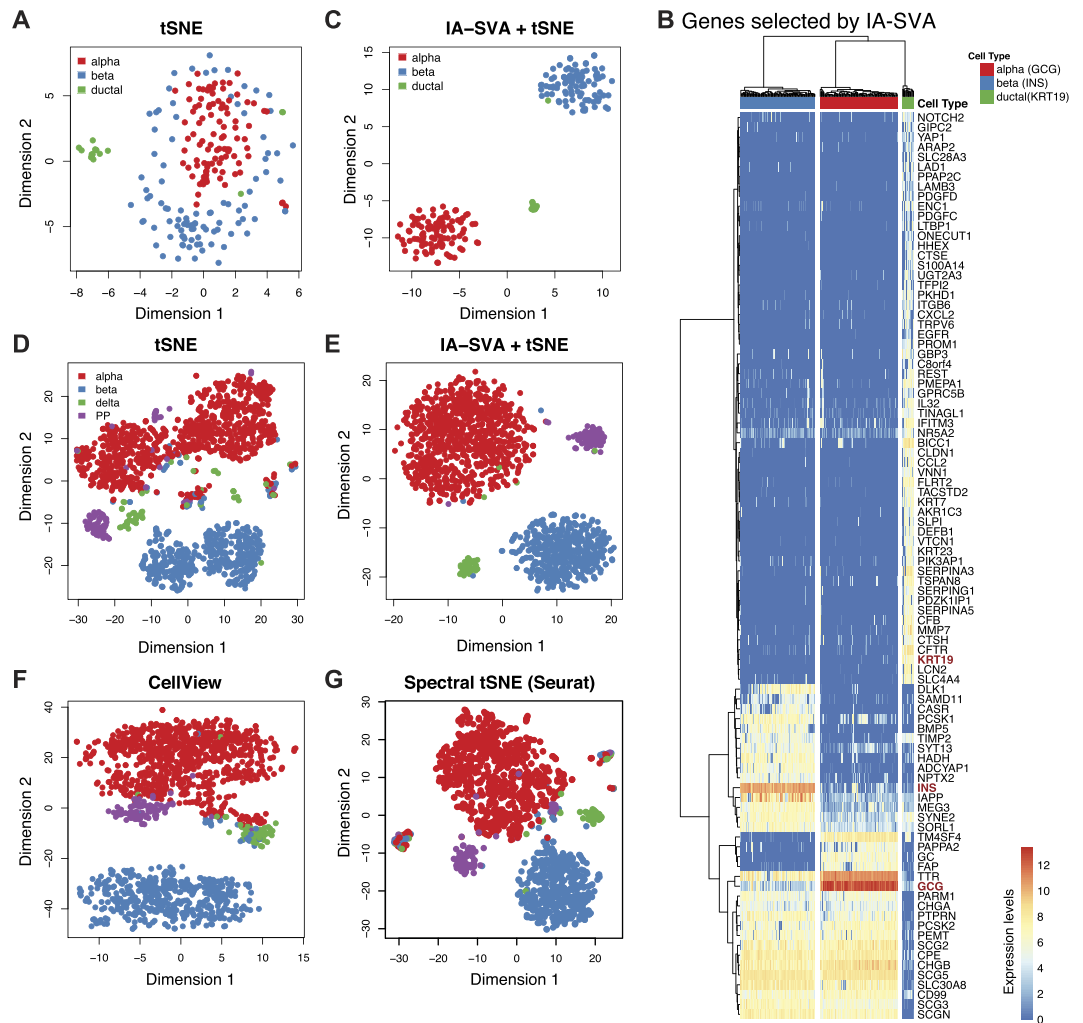


Figure 4. IA-SVA based gene selection enhances single cell data visualization. **(A)** tSNE analyses using all expressed genes ($n = 19,226$) in human islet data (tSNE). Cells are color-coded based on the original cell-type assignments. Note that cells are not effectively clustered with respect to their assigned cell types. **(B)** Hierarchical clustering using genes ($n = 92$) selected by IA-SVA clearly separate cell types (ward.D2 and `cutree_cols = 3`). Known marker genes (e.g., *INS*) are highlighted in red color. The values in the heatmap and the color bar are log-transformed (base e) normalized counts (normalized by dividing each cell by its total counts then multiplying median of library size). **(C)** tSNE analyses using the 92 IA-SVA genes (IA-SVA + tSNE). Note the improved grouping of cell types into discrete clusters. **(D)** tSNE analyses using top variable genes in a second and larger islet scRNA-seq dataset. Note that cells are not effectively clustered with respect to their assigned cell types just using tSNE. **(E)** tSNE analyses repeated using genes ($n = 57$) obtained via IA-SVA (IA-SVA + tSNE). Note the improved clustering of different cell types into discrete clusters. **(F)** tSNE analyses using 1000 most over-dispersed genes (CellView). **(G)** tSNE analyses on significant PCs obtained from highly over-dispersed genes (Spectral tSNE).

Testing the performance of IA-SVA and existing methods when known and hidden factors affect overlapping sets of genes.

To study the impact of a significant overlap between hidden and known factors on IA-SVA's performance, we performed an analysis, where we simulated gene expression matrices with 10,000 genes and 50 cells using the Polyester R package in addition to a known factor affecting 2,000 genes and a hidden factor affecting 1,000 genes, where these two factors are highly correlated ($r = 0.76$). The effect size of genes affected by each factor is randomly simulated from a standard normal distribution as in the previous simulations. Next, we studied IA-SVA's performance for inferring the hidden factor when we change the percent of genes affected by the hidden factor that are also affected by the known factor, ($k\%$, $k \in \{99, 90, 80, 70, 60, 50, 40, 30, 20, 10, 0\}$), where $k = 99\%$ is the most extreme case where 99% of all genes affected by the hidden factor are also affected by the known factor. Under each setting, we estimated the hidden factor using IA-SVA, PCA, USVA and SVA and computed the absolute correlation coefficients between the estimates and the true hidden factor. As shown in the Fig. S9, IA-SVA accurately captures the simulated factor even under the extreme scenario where 99% of genes affected by the hidden factor are also affected by the known factor. The accuracy of USVA and SVA declines as the overlap percentage goes up. These analyses suggest that even in extreme scenarios where known

and hidden factors affect overlapping gene sets IA-SVA can detect the hidden heterogeneity with high accuracy. We further studied the data generated under the most extreme scenario (i.e., $k = 99\%$) to understand what leads to IA-SVA's increased performance and noted that even if the same set of genes are affected by both known and hidden factors, the effect sizes can vary, where for some genes effect sizes are higher for the known factor and for some others effect sizes are higher for hidden factors (Fig. S10). Among the genes affected by both factors, genes strongly affected by the hidden factor yet weakly affected by the known factor tends to have higher R^2 values (Fig. S10D,E). This heterogeneity in effect sizes enables IA-SVA to capture genes strongly associated with the hidden factor and estimate it with high accuracy.

Time complexity analyses for alternative IA-SVA implementations. We measured the run time of IA-SVA and its alternative implementations for increasing number of genes and cells in the input data using both serial (“threads” parameter in the R package = 1) and parallel settings (threads > 1). We first run IA-SVA in serial mode using 4 different IA-SVA modes: 1) Default (i.e., full SVD, number of permutations (num.p) = 100); 2) Partial-SVD (uses 10 (i.e., num.sv.permtest = 10) largest singular values/vectors obtained by using the augmented implicitly restarted Lanczos bidiagonalization algorithm (IRLBA) (<https://bwlewis.github.io/irlba/>) in permutation-based significance index (SI) calculations, num.p = 100); 3) Fast (fast_iasva function), that keeps hidden factors explaining at least $k\%$ of total variance (default k (pct.cutoff) = 1%) instead of the permutation based significance test; and 4) No Permutation (permute = FALSE) for each dataset. Each run was repeated ten times. As expected IA-SVA runs faster without the permutations (Supplementary Fig. S11A “Fast” and “No permute” settings) and takes less than a few minutes to finish even for large datasets composed of thousands of cells using these two settings. Furthermore, these results suggest that IA-SVA's run time is approximately quadratic for “full SVD” and “fast IA-SVA” modes in terms of cell numbers ($O \sim m^2$, where m = the number of cells) and linear for “No permutation” and “Partial SVD” modes ($O \sim m$) (Supplementary Fig. S11B). On the other hand, its run time is linear for all four modes with increasing number of genes ($O \sim n$, where n = the number of genes) (plot not shown).

We implemented the parallelization option in the IA-SVA function (“threads” parameter) for full and partial SVD settings to reduce the run time for permutation based significance testing. Here we tested how parallelizable IA-SVA is by running full and partial SVD modes using different number of cores (1 to 16) on two different datasets: a gene-dense dataset (16000 genes and 640 cells) and a cell-dense dataset (1000 genes and 5120 cells). For each IA-SVA run, we used the geometric library size and patient identifier vector as known covariates in the model. To be able to compare alternative implementations of IA-SVA, we only detected the first SV (num.sv = 1). To benchmark single-cell analysis in a high performance-computing environment, multiple IA-SVA runs (parameter sweeps across the number of genes and cells) were submitted as a single job (i.e., analysis group). For each IA-SVA implementation, each analysis group was repeated ten times by submitting to nodes with minor variability in their cluster environments (HP Proliant SL/XL Series servers with 2.5–2.6 GHz processor clock speed and 256GB RAM). To capture the variation stemming from the heterogeneous cluster environment, the standard error of mean for each parameter (gene or cell) was reported. Our analyses showed that parallelization significantly reduced IA-SVA's run time (5 to 8-fold speed up with 16 cores), where Amdahl's law suggests that approximately 90–95% of runtime is parallelizable (Supplementary Fig. S11C). These analyses suggest that the most time-consuming step in IA-SVA is the calculation of SI values. Therefore, for larger datasets (i.e., 10X genomics), we recommend running the job on multiple cores or using the fast IA-SVA function (fast_iasva), which selects important SVs based on the proportion of variation explained. Under these circumstances, the run time of IA-SVA should not be a limiting factor in data analyses.

Discussion

Surrogate variable analyses (SVA) based methods are effective in detecting and eliminating hidden and unwanted variation in bulk gene expression data (such as batch effects). These methods are not developed to accurately estimate the hidden factors; instead they assume that hidden factors are typically unwanted (e.g., batch effects) and they aim to effectively adjust for the unwanted variation without compromising the biological signal of interest (e.g., group variable: case vs control) in downstream analyses (e.g., differential analyses)¹⁴. Although these methods are very effective in batch correction; measurement of gene expression levels at single cell resolution pose novel challenges in the detection and adjustment of hidden sources of data variation. First, single cell transcriptomes harbour hidden variation that can be biologically interesting (hence ‘wanted’) and therefore, accurately estimating hidden ‘wanted’ factor(s) can be the major goal of a study, for example detection of rare cells within a tissue²⁹ or detection of cell subtypes that can be linked to health or disease²⁸. Second, since single cell data do not average out variation as in the case of bulk profiling, the data reflect variation arising from diverse biological and technical sources some of which could be highly correlated. Existing SVA-based methods do not readily apply to the unique needs of single cell data analyses.

To fill this gap, we developed IA-SVA, which is designed to accurately estimate hidden and biologically interesting (i.e., ‘wanted’) factors even if these factors are correlated with each other and/or with the known factors. Unlike other SVA-based methods, IA-SVA therefore focuses more on the accurate detection and estimation of hidden factors rather than their elimination since these factors can be biologically interesting, e.g., identification of a new cell type and its marker genes. Indeed, analyses on simulated scRNA-seq data showed that IA-SVA outperforms existing supervised (i.e., SSVA) and unsupervised (i.e., USVA) state-of-the-art methods in the estimation of hidden factors (not necessarily in their elimination). Furthermore, we noted that IA-SVA is particularly effective (i.e., high detection power and accuracy, and Type I error rate controlled under the nominal level of 0.05) in detecting correlated factors that affect a small fraction of genes. Therefore IA-SVA is an effective unsupervised alternative to existing SVA-based algorithms when the goal is to accurately estimate hidden factors (and

their marker genes) rather than to eliminate these factors — if the goal of an analysis is to identify and strip the data from any batch effects (hidden and unwanted variation) while preserving the variation associated with the primary variable of interest, existing algorithms (e.g., SVA) can robustly and efficiently accomplish this task and IA-SVA should not be the method of choice.

Through analyses of diverse human datasets from multiple studies, we established that IA-SVA can effectively detect hidden heterogeneity in scRNA-seq data arising from a small number of cells either due to technical (i.e., contamination or doublets) or biological (i.e., a rare cell type) sources. In two independently generated islet scRNA-seq datasets, we showed that IA-SVA detects heterogeneity stemming from compromised alpha cells (which may be explained by cell contamination), which should be excluded from the downstream analyses (Fig. 2 and Supplementary Fig. S3). Therefore, IA-SVA provides an easy-to-apply statistical framework to uncover variation in scRNA-seq data even if it is stemming from only a handful of cells. This ability of IA-SVA can be effective in identifying rare cells within a population of cells, where genes associated with the detected factor can uncover relevant marker genes for the rare population of cells. In addition, IA-SVA can be effective in detecting heterogeneity associated with cell-cycle stages without prior knowledge, therefore providing an unsupervised solution to this common problem in single cell data analyses (Fig. 3).

An important feature of IA-SVA is its ability to uncover genes associated with detected hidden factors. This feature can be used to detect marker genes associated with different cell types. As a proof-of-concept we demonstrated this in pancreatic islet cells, where we captured known marker genes (e.g., *INS*, *GCG*) in an unsupervised manner. Moreover, genes captured by IA-SVA can be used to improve the visualization of single cells into their respective clusters, as demonstrated with the analyses of islet cells from two separate studies (Fig. 4). Spectral tSNE¹⁶ is a commonly used method for scRNA-seq data visualization especially in the existence of confounding factors. This method regresses out variation associated with known factors before data visualization. However, when a hidden factor is ‘wanted’ (e.g., cell types) and is highly correlated with known factors, removing the known factors will also diminish the ability to detect the wanted hidden factor and the genes associated with this factor (e.g., marker genes for different cell types). Indeed, our analyses using islet cells emphasized the importance of properly adjusting the data for known factors prior to further analyses, such as data visualization (e.g., tSNE) to prevent spurious clustering of cells due to the confounding factors (Fig. 4E). IA-SVA is an alternative method that can effectively handle data with multiple confounding factors.

One caveat of IA-SVA analyses arises from its permutation-based quasi-inferential procedure to determine the number of hidden factors similar to the permutation-based parallel analysis³⁰. These measurements, referred to as significance index (SI) rather than empirical p-values since the p-values obtained from the quasi-inferential procedure are only meaningful for the largest eigenvalue and it is not straightforward to adjust them for multiple hypothesis testing. We recommend users of IA-SVA to select a stringent cut-off for SI values especially if they expect to detect many unknown factors from their data. Our simulation and real-world data analyses suggest that despite its caveats SI is still useful in determining the number of hidden factors in single cell data analyses.

In summary, IA-SVA is an SVA-based unsupervised method designed to accurately estimate hidden factors (sources of variation) in single cell gene expression data. The iterative and flexible framework of IA-SVA allows the accurate estimation of multiple and potentially correlated factors along with their statistical significance, which is the main advantage of IA-SVA over existing methods. This flexibility is more realistic given the confounded nature of known and unknown factors in single cell gene expression measurements. Therefore, IA-SVA has an improved performance over existing SVA-based methods in terms of estimating hidden sources of variation when they are correlated with each other and with known variables. IA-SVA is also an effective alternative to methods developed for single cell data analyses (e.g., CellView, ZINB-WaVe, and Seurat), especially for the analyses of complex data (i.e., data with multiple confounding and correlated factors). With the increasing amount of single cell studies and the increasing complexity of human cohorts, IA-SVA will serve as an effective statistical framework specifically designed to handle unique challenges of scRNA-seq data analyses. Computational method development for single cell datasets is a very active research area, where different methods have been developed to extract and interpret heterogeneity from such data, e.g., methods based on non-negative matrix factorization^{31,32}. Effectively comparing all of these methods will require building benchmark datasets and implementing these methods using similar platforms. There is an ongoing consortium effort by the Human Cell Atlas project (<https://www.humancellatlas.org/>) to systemically benchmark computational methods developed for single cell data analyses, including IA-SVA.

Methods

IA-SVA framework. We model the log-transformed sequencing read counts for m cells and n genes (i.e., $Y_{m \times n}$) as a combination of known and unknown variables as follows:

$$Y_{m \times n} = X_{m \times p} \beta_{p \times n} + Z_{m \times k} \delta_{k \times n} + \varepsilon_{m \times n},$$

where $X_{m \times p}$ is a matrix for p known variable(s) (e.g., patient ID, sex or ethnicity), $Z_{m \times k}$ is a matrix for k unknown variables and $\varepsilon_{m \times n}$ is the error term. With this model, we can account for any clinical/experimental information about samples (e.g., sex, ethnicity, age, BMI or batch) as known factors ($X_{m \times p}$) and dissect unaccounted variation in the read count data that is attributable to hidden factors ($Z_{m \times k}$).

Existing unsupervised SVA-based methods (e.g., USVA¹⁰, RUV¹³, ISVA¹²) obtain the residual matrix ($Y'_{m \times n}$) by regressing log-transformed read counts ($Y_{m \times n}$) on all known factors ($X_{m \times p}$). Then, they infer the hidden factors from this residual matrix ($Y'_{m \times n}$) using dimensionality reduction algorithms (e.g., PCA, SVD or ICA). Thus, by definition, multiple hidden factors captured by these methods are orthogonal to each other and to known variables. Therefore, if hidden factors are correlated with each other and with known factors, the direct inference from the residual matrix leads to biased estimates of hidden factors due to the orthogonality assumption.

In contrast, IA-SVA does not impose orthogonality between factors (hidden or known) and estimates correlated factors via a novel iterative framework (Fig. 1). At each iteration, IA-SVA first obtains residual matrix ($Y'_{m \times n}$), i.e., log-transformed read count matrix adjusted for all known factors ($X_{m \times p}$) including surrogate variables of unknown factors estimated from previous iterations and extracts the first principal component (PC1) from the residuals ($Y'_{m \times n}$) using SVD. Next it tests the significance of PC1 in terms of its contribution to the unmodeled variation (i.e., residual variance). Using this PC1 as a surrogate variable (as in the case of existing methods) implicitly imposes orthogonality between known and hidden factors. Instead, IA-SVA uses PC1 to infer gene weights, which are also used to infer genes associated with the hidden factor. IA-SVA relies on the fact that the first principal component (PC1) of the residual matrix is highly correlated with the hidden factor that contributes the most to the unmodeled variation in the data, and thus, PC1 can be used to sort genes in terms of their relative association strength with the hidden factor. To infer these genes, IA-SVA regresses the residual matrix Y' on PC1 and calculates the coefficient of determination (R^2) for each gene. Genes with high R^2 scores can be treated as marker genes for the factor. These R^2 scores are further utilized for an unbiased inference of the hidden factor while retaining the correlation structure between known and hidden factors. For this, IA-SVA first obtains a weighted log-count matrix ($Y''_{m \times n}$) by weighing all genes with respect to their R^2 scores (i.e., $Y''_{m \times n} = Y_{m \times n} W_{n \times n}$, where Y is the original log-transformed read counts and W is a diagonal matrix of R^2 values). Then it conducts a SVD on the weighted matrix $Y''_{m \times n}$ and obtains the PC1 to be used as a surrogate variable (SV) for the hidden factor. In the next iteration, IA-SVA uses this SV as an additional known factor to identify further significant hidden factors.

The rationale for using the coefficient of determination (R^2) in our weighting scheme is two-folded. First, R^2 indicates the proportion of the variance in gene read counts (dependent variable) that can be explained by the PC1 obtained from the residual matrix in Step 2 (independent variable). Therefore, R^2 is useful to infer genes whose gene expression variation is most strongly explained by PC1, hence by the hidden factor, based on the assumption that PC1 strongly correlates with the hidden factor. Therefore, R^2 is an intuitive choice to rank genes based on their relevance to the hidden factor. Second, R^2 ranges from 0 to 1, which would be appropriate to use as weights, where higher values imply higher relevance to PC1, hence the hidden factor. Any other measure that would meet these criteria would also work in our methodology. Instead of using the weighting approach, one could also conduct association testing by regressing the residual matrix on PC1, identify significant genes (e.g., FDR < 0.05) with high effect sizes (e.g., $R^2 > 0.3$), conduct a SVD on the read counts of these significant genes to obtain PC1 and finally use this PC1 as the hidden factor estimate. In our simulation studies, we experimented with this alternative schema and observed that R^2 -based weighting provides more robust results across different settings (data not shown). We further illustrate with more details how the R^2 -based weighting schema works in Supplementary Text S1.

The iterative procedure of IA-SVA is composed of six major steps as summarized in Fig. 1A and below:

[Step 1] Regress $Y_{m \times n}$ on all known factors ($X_{m \times p}$), including SVs obtained from previous iterations, to obtain residuals ($Y'_{m \times n}$).

[Step 2] Conduct SVD on the obtained residuals ($Y'_{m \times n}$) and obtain the first PC (PC1).

[Step 3] Test the significance of the contribution of PC1 to the variation in residuals ($Y'_{m \times n}$) using a non-parametric permutation-based assessment^{9,10,30} as explained further in the next section.

[Step 4] If PC1 is significant, regress the residual matrix $Y'_{m \times n}$ on PC1 to compute the coefficient of determination (R^2) for every gene. If PC1 is not significant, stop the iteration and conduct subsequent downstream analysis using previously obtained significant SVs.

[Step 5] Weigh each gene in $Y_{m \times n}$ with respect to its R^2 value by multiplying a gene's log-transformed read counts ($Y_{m \times n}$) with its R^2 values ($Y''_{m \times n} = Y_{m \times n} W_{n \times n}$). The highly weighted genes in this framework serve as the genes affected by the hidden factor.

[Step 6] Conduct a second SVD on this weighted log-counts matrix ($Y''_{m \times n} = Y_{m \times n} W_{n \times n}$) to obtain PC1, which will be used as the surrogate variable (SV) for the hidden factor.

At the end of this six-step procedure, IA-SVA uses the detected SV (if significant) as an additional known factor in the next iteration. The algorithm stops, when no more significant PC1s are detected in Step 3. Significant SVs obtained via IA-SVA can be used in subsequent analyses. If an SV arises from an unwanted factor (e.g., cell contamination), these SVs can be included as covariates in the model to remove the unwanted variation or to filter out contaminated cells. In single cell data significant SVs could also explain 'wanted' biological factors (e.g., different cell types) and genes associated with such SVs can be further evaluated to discover novel biology from these complex datasets.

Quasi-inferential procedure for determining the number of hidden factors. We implemented a permutation based quasi-inferential procedure to determine the number of SVs to be considered in downstream data analyses and interpretation. For this we implemented a permutation based method that empirically assess the significance of the contribution of a hidden factor estimate (i.e., PC1 obtained in Step 2) to the residual variation — calculating the proportion of variation explained by the PC1 and conducting a permutation based significance test similar to parallel analysis³⁰ and surrogate variable analysis¹⁰. Unlike parallel analysis and SVA¹⁰, which tests all putative hidden factors at once, IA-SVA assesses the significance of hidden factors one at a time during the corresponding iteration (always for the PC1 of the residual matrix $Y'_{m \times n}$ detected in that iteration). Briefly, IA-SVA i) conducts SVD on the residual matrix obtained from Step 1, ii) computes the proportion of variation in this matrix explained by the first singular vector (i.e., PC1) and iii) compares this proportion against the values obtained from permuted residual matrices, as further explained below:

[Step 1] Conduct SVD on the residual matrix ($Y'_{m \times n}$).

[Step 2] Calculate the proportion of residual variance explained by the first singular vector (PC1) using the test statistic: $T_{obs} = \frac{\lambda_1^2}{\sum_k \lambda_k^2}$, where λ_k is the k -th singular value.

[Step 3] Generate a permuted residual matrix by i) permuting each column of the log-transformed read count matrix $Y_{m \times n}$ and regressing the permuted read count matrix on all known factors ($X_{m \times p}$) to obtain fitted residuals.

[Step 4] Repeat Step 3 M times and generate an empirical null distribution of the test statistic by calculating (T_i^0 , $i = 1, \dots, M$) for the M permuted residual matrices.

[Step 5] Compute a proportion (named as significance index) for the first singular vector (PC1) by counting the number of times the null statistics (T_i^0) exceeds the observed one (T_{obs}) divided by the number of permutations (M).

We prefer to call the proportion obtained in Step 5 “significance index”, not “empirical p-value” since the quasi-inferential procedure does not perfectly meet the requirements of the classical hypothesis testing similar to the permutation based parallel analysis³⁰ (e.g., the permutation p-value is only statistically meaningful for the largest eigenvalue).

scRNA-seq data simulation. To eliminate the potential bias in data simulations and make simulation studies more objective³³, we used a third-party simulation software (i.e., Polyester R package³⁴) and study design (<http://jtleek.com/svaseq/simulateData.html>) and simulated scRNA-seq data to test IA-SVA's performance. The original simulation design is slightly modified to reflect characteristics of scRNA-seq data for high dropout rate (i.e., excessive number of zeros in the data), high variability (i.e., over-dispersed gene counts) and multiple hidden factors highly correlated with known factors. First, to simulate high dropout rates (proportion of zero counts $\approx 70\%$) and high variability, we estimated Polyester's zero-inflated negative binomial model parameters (i.e., p_0 : probabilities that the count will be zero, μ : mean of the negative binomial, $size$: size of the negative binomial) from human pancreatic islet data generated using the Fluidigm C1 platform¹⁹. Using these estimated model parameters, we simulated expression data for m cells and n genes under two hypotheses: (1) the null hypothesis: no hidden sources of variation, and (2) the alternative hypothesis: three hidden factors with two levels (e.g., cell contamination factor: contaminated vs. normal cells). For both scenarios, baseline expression levels were simulated using the Polyester R package by using the zero-inflated negative binomial distribution parameters estimated from the islet scRNA-seq data. Under both scenarios, we simulated a known factor with two levels (e.g., males vs. females) and simulated 10% of genes to be differentially expressed between the two levels. Under the alternative hypothesis, we simulated three two-level hidden factors that affect 30%, 20% and 10% of randomly chosen genes respectively and simulated two different scenarios where these factors are moderately correlated ($|r| = \sim 0.3-0.6$) or weakly correlated ($|r| < 0.3$) with the known variable. For genes affected by a factor (both for known and hidden factors), effect sizes (regression coefficients, i.e., β or δ) of the factor were randomly simulated from a standard normal distribution and the effect sizes for un-affected genes were set to zero¹¹. The R code for data simulations is publicly available for reproducibility at https://dleelab.github.io/iasvaExamples/scRNAseq_simulation.html.

Data processing and normalization. In all analyses, we used genes with >5 reads in at least 3 cells (i.e., a gene has to be above five in at least three cells to be kept) and normalized the retained gene expression counts using SCnorm³⁵ with default settings for further analyses. For single cell data visualization examples, we normalized gene read counts by dividing each cell column by its total counts then multiplying median of library size, which is similar to the default normalization method (“LogNormalize”) implemented in Seurat¹⁶.

Data simulations for time complexity analyses. Pancreatic islet data¹⁹ that is composed of 18,177 genes and 638 cells was replicated ten times to increase cell numbers and a fixed seed was used to randomly down-sample this data into smaller matrix of n genes and m cells, where $n \in \{1000, 2000, 4000, 8000\}$ and $m \in \{20, 40, 80, 160, 320, 640, 1280, 2560, 5120\}$. For studying the parallelization of IA-SVA, pancreatic islet data was down-sampled into a gene-dense (composed of 16,000 genes and 640 cells) and a cell-dense (composed of 1000 genes and 5120 cells) data.

Availability of Data and Methods

An R package for IA-SVA with example case scenarios is freely available from <https://github.com/UcarLab/iasva>. This R package has been accepted by Bioconductor and is currently available at <https://www.bioconductor.org/packages/devel/bioc/html/iasva.html>. The published data sets analyzed in this study including single-cell RNA sequencing read counts and annotations describing samples and experiment settings are included in an R data package (iasvaExamples) deposited at <https://github.com/dleelab/iasvaExamples>. We also implemented an R Shiny web application that conducts gene enrichment analyses using publicly available databases (e.g., GO terms, KEGG pathways, immune modules, cell-type specific genes) on a set of genes significantly associated with each SV (https://nlawlor.shinyapps.io/IASVA_Shiny_08_13_2018/). This application also includes methods for convenient clustering and visualization of data. Our tool can also be used with other factor analyses to help interpret their outputs.

References

1. Tung, P.-Y. *et al.* Batch effects and the effective design of single-cell gene expression studies. *bioRxiv*, 062919 (2016).
2. Kowalczyk, M. S. *et al.* Single-cell RNA-seq reveals changes in cell cycle and differentiation programs upon aging of hematopoietic stem cells. *Genome research* **25**, 1860–1872 (2015).
3. Stegle, O., Teichmann, S. A. & Marioni, J. C. Computational and analytical challenges in single-cell transcriptomics. *Nat Rev Genet* **16**, 133–145, <https://doi.org/10.1038/nrg3833> (2015).
4. Villani, A. C. *et al.* Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors. *Science* **356**, <https://doi.org/10.1126/science.aah4573> (2017).
5. Ilicic, T. *et al.* Classification of low quality cells from single-cell RNA-seq data. *Genome Biol* **17**, 29, <https://doi.org/10.1186/s13059-016-0888-1> (2016).
6. Buettner, F. *et al.* Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol* **33**, 155–160, <https://doi.org/10.1038/nbt.3102> (2015).
7. McDavid, A., Finak, G. & Gottardo, R. The contribution of cell cycle to heterogeneity in single-cell RNA-seq data. *Nat Biotechnol* **34**, 591–593, <https://doi.org/10.1038/nbt.3498> (2016).
8. Hicks, S. C., Teng, M. & Irizarry, R. A. On the widespread and critical impact of systematic bias and batch effects in single-cell RNA-Seq data. *bioRxiv*, <https://doi.org/10.1101/025528> (2015).
9. Leek, J. T. & Storey, J. D. Capturing heterogeneity in gene expression studies by surrogate variable analysis. *PLoS Genet* **3**, 1724–1735, <https://doi.org/10.1371/journal.pgen.0030161> (2007).
10. Leek, J. T. & Storey, J. D. A general framework for multiple testing dependence. *Proc Natl Acad Sci USA* **105**, 18718–18723, <https://doi.org/10.1073/pnas.0808709105> (2008).
11. Leek, J. T. svaseq: removing batch effects and other unwanted noise from sequencing data. *Nucleic Acids Res* **42**, <https://doi.org/10.1093/nar/gku864> (2014).
12. Teschendorff, A. E., Zhuang, J. & Widschwendter, M. Independent surrogate variable analysis to deconvolve confounding factors in large-scale microarray profiling studies. *Bioinformatics* **27**, 1496–1505, <https://doi.org/10.1093/bioinformatics/btr171> (2011).
13. Risso, D., Ngai, J., Speed, T. P. & Dudoit, S. Normalization of RNA-seq data using factor analysis of control genes or samples. *Nat Biotechnol* **32**, 896–902, <https://doi.org/10.1038/nbt.2931> (2014).
14. Gagnon-Bartsch, J. A. & Speed, T. P. Using control genes to correct for unwanted variation in microarray data. *Biostatistics* **13**, 539–552, <https://doi.org/10.1093/biostatistics/kxr034> (2012).
15. Maaten, L. V. D. Accelerating t-SNE using tree-based algorithms. *J. Mach. Learn. Res.* **15**, 3221–3245 (2014).
16. Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161**, 1202–1214, <https://doi.org/10.1016/j.cell.2015.05.002> (2015).
17. Risso, D., Perraudeau, F., Gribkova, S., Dudoit, S. & Vert, J. P. A general and flexible method for signal extraction from single-cell RNA-seq data. *Nat Commun* **9**, 284, <https://doi.org/10.1038/s41467-017-02554-5> (2018).
18. Darmanis, S. *et al.* A survey of human brain transcriptome diversity at the single cell level. *Proc Natl Acad Sci USA* **112**, 7285–7290, <https://doi.org/10.1073/pnas.1507125112> (2015).
19. Lawlor, N. *et al.* Single cell transcriptomes identify human islet cell signatures and reveal cell-type-specific expression changes in type 2 diabetes. *Genome Res*, <https://doi.org/10.1101/gr.212720.116> (2016).
20. Xin, Y. *et al.* RNA Sequencing of Single Human Islet Cells Reveals Type 2 Diabetes Genes. *Cell Metab* **24**, 608–615, <https://doi.org/10.1016/j.cmet.2016.08.018> (2016).
21. Xin, Y. *et al.* Use of the Fluidigm C1 platform for RNA sequencing of single mouse pancreatic islet cells. *Proc Natl Acad Sci USA* **113**, 3293–3298, <https://doi.org/10.1073/pnas.1602306113> (2016).
22. Wang, Y. J. *et al.* Single-Cell Transcriptomics of the Human Endocrine Pancreas. *Diabetes* **65**, 3028–3038, <https://doi.org/10.2337/db16-0405> (2016).
23. Patel, A. P. *et al.* Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* **344**, 1396–1401, <https://doi.org/10.1126/science.1254257> (2014).
24. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Res* **44**, D457–462, <https://doi.org/10.1093/nar/gkv1070> (2016).
25. Gene Ontology, C. Gene Ontology Consortium: going forward. *Nucleic Acids Res* **43**, D1049–1056, <https://doi.org/10.1093/nar/gku1179> (2015).
26. Lun, A. T., McCarthy, D. J. & Marioni, J. C. A step-by-step workflow for low-level analysis of single-cell RNA-seq data with Bioconductor. *F1000Res* **5**, 2122, <https://doi.org/10.12688/f1000research.9501.2> (2016).
27. Bolisetty, M. T., Stitzel, M. L. & Robson, P. CellView: Interactive Exploration Of High Dimensional Single Cell RNA-SeqData. *bioRxiv*, <https://doi.org/10.1101/123810> (2017).
28. Lawlor, N., Khetan, S., Ucar, D. & Stitzel, M. L. Genomics of Islet (Dys)function and Type 2 Diabetes. *Trends Genet* **33**, 244–255, <https://doi.org/10.1016/j.tig.2017.01.010> (2017).
29. Prosperio, V. & Lonngberg, T. Single-cell technologies are revolutionizing the approach to rare cells. *Immunol Cell Biol* **94**, 225–229, <https://doi.org/10.1038/icb.2015.106> (2016).
30. Buja, A. & Eyuboglu, N. Remarks on Parallel Analysis. *Multivariate Behav Res* **27**, 509–540, https://doi.org/10.1207/s15327906mbr2704_2 (1992).
31. Zhu, X., Ching, T., Pan, X., Weissman, S. M. & Garmire, L. Detecting heterogeneity in single-cell RNA-Seq data by non-negative matrix factorization. *PeerJ* **5**, e2888, <https://doi.org/10.7717/peerj.2888> (2017).
32. Stein-O'Brien, G. L. *et al.* Decomposing cell identity for transfer learning across cellular measurements, platforms, tissues, and species. *bioRxiv*, <https://doi.org/10.1101/395004> (2018).
33. Gelman, A. & Hennig, C. Beyond subjective and objective in statistics. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, n/a–n/a, <https://doi.org/10.1111/rssa.12276> (2017).
34. Frazee, A. C., Jaffe, A. E., Langmead, B. & Leek, J. T. Polyester: simulating RNA-seq datasets with differential transcript expression. *Bioinformatics* **31**, 2778–2784, <https://doi.org/10.1093/bioinformatics/btv272> (2015).
35. Bacher, R. *et al.* SCnorm: robust normalization of single-cell RNA-seq data. *Nat Methods* **14**, 584–586, <https://doi.org/10.1038/nmeth.4263> (2017).

Acknowledgements

We thank the Jackson Laboratory Computational Science group, Ucar and Stitzel lab members for constructive feedback throughout this project. We thank Jane Cha, JAX scientific illustrator, for her help with Figure 1. This study was made possible by generous financial support of the National Institute of General Medical Sciences (NIGMS) under award number GM124922 (to D.U.) and of the Chan Zuckerberg Initiative (CZI) under award number 182753 SVCF (to D.U.) and by the Jackson Laboratory Scientific Services Innovation Fund (to D.L. and D.U.) Opinions, interpretations, conclusions, and recommendations are solely the responsibility of the authors and do not necessarily represent the official views of the National Institutes of Health (NIH).

Author Contributions

D.L. and D.U. designed the project, generated the figures and wrote the manuscript. D.L. developed the statistical framework and run the data analyses. M.B. provided advice on data analyses and interpretation of results. A.C. contributed to the data pre-processing and the generation of the R package. N.L. formatted and submitted the R package to Bioconductor. All authors read and approved this manuscript prior to submission.

Additional Information

Supplementary information accompanies this paper at <https://doi.org/10.1038/s41598-018-35365-9>.

Competing Interests: The authors declare no competing interests.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2018