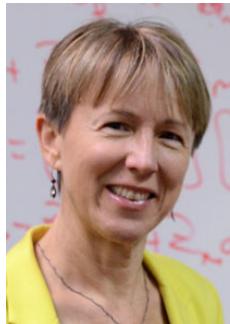


Carnegie Mellon University
Statistics & Data Science

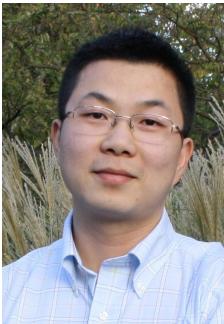
High-dimensional statistical methods to model heterogeneity in genomic data

Kevin Lin

Advisors



Kathryn
Roeder



Jing
Lei



Max
G'Sell

Committee Members



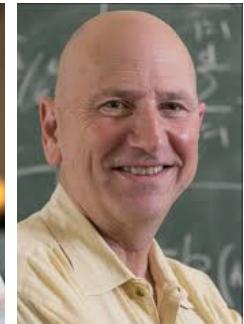
Han
Liu



Anjali
Mazumder



Ryan J.
Tibshirani



Larry
Wasserman

Carnegie Mellon University
Statistics & Data Science

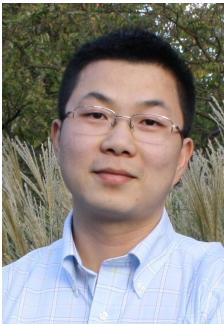
High-dimensional statistical methods to model heterogeneity in genomic data

Kevin Lin

Advisors



Kathryn
Roeder



Jing
Lei



Max
G'Sell

Committee Members



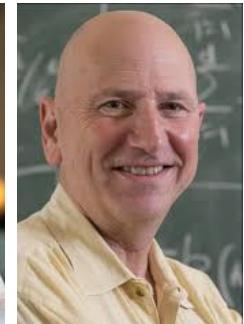
Han
Liu



Anjali
Mazumder



Ryan J.
Tibshirani



Larry
Wasserman

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 2: Studying risk genes for autism, by finding a subset of homogeneous samples within a heterogeneous dataset via multiple covariance testing

(Lin, Liu, Roeder: JASA, 2020)

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 2: Studying risk genes for autism, by finding a subset of homogeneous samples within a heterogeneous dataset via multiple covariance testing

(Lin, Liu, Roeder: JASA, 2020)

Chapter 3-4: Studying copy number variation, by proving properties of changepoint estimation and inference

(Lin, Sharpnack, Rinaldo, Tibshriani: NeurIPS, 2017)

(Hyun, Lin, G'Sell, Tibshriani: Biometrics, In revision)

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 2:

Studying risk genes for autism, by finding a subset of homogeneous samples within a heterogeneous dataset via multiple covariance testing

(Lin, Liu, Roeder: JASA, 2020)

Chapter 3-4:

Studying copy number variation, by proving properties of changepoint estimation and inference

(Lin, Sharpnack, Rinaldo, Tibshirani: NeurIPS, 2017)

(Hyun, Lin, G'Sell, Tibshirani: Biometrics, In revision)

Chapter 5:

Studying developmental trajectories of oligodendrocytes, by developing non-linear embedding based on exponential-family distributions

(Lin, Lei, Roeder: JASA, In revision)

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 5:

Studying developmental trajectories of oligodendrocytes, by developing non-linear embedding based on exponential-family distributions

(Lin, Lei, Roeder: JASA, In revision)

Focus of this talk

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 5:

Studying developmental trajectories of oligodendrocytes, by developing non-linear embedding based on exponential-family distributions

(Lin, Lei, Roeder: JASA, In revision)

Focus of this talk

Punch line (Statistics): eSVD performs a non-linear embedding by estimating a low-rank matrix of natural parameters via alternating minimization based on a random dot product model.

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 5:

Studying developmental trajectories of oligodendrocytes, by developing non-linear embedding based on exponential-family distributions

(Lin, Lei, Roeder: JASA, In revision)

Focus of this talk

Punch line (Statistics): eSVD performs a non-linear embedding by estimating a low-rank matrix of natural parameters via alternating minimization based on a random dot product model.

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 5:

Studying developmental trajectories of oligodendrocytes, by developing non-linear embedding based on exponential-family distributions

(Lin, Lei, Roeder: JASA, In revision)

Focus of this talk

Punch line (Statistics): eSVD performs a non-linear embedding by estimating a low-rank matrix of natural parameters via alternating minimization based on a random dot product model.

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 5:

Studying developmental trajectories of oligodendrocytes, by developing non-linear embedding based on exponential-family distributions

(Lin, Lei, Roeder: JASA, In revision)

Focus of this talk

Punch line (Statistics): eSVD performs a non-linear embedding by estimating a low-rank matrix of natural parameters via alternating minimization based on a random dot product model.

Biologists model the heterogeneity in genomic data in order to investigate different biological processes.

Chapter 5:

Studying developmental trajectories of oligodendrocytes, by developing non-linear embedding based on exponential-family distributions

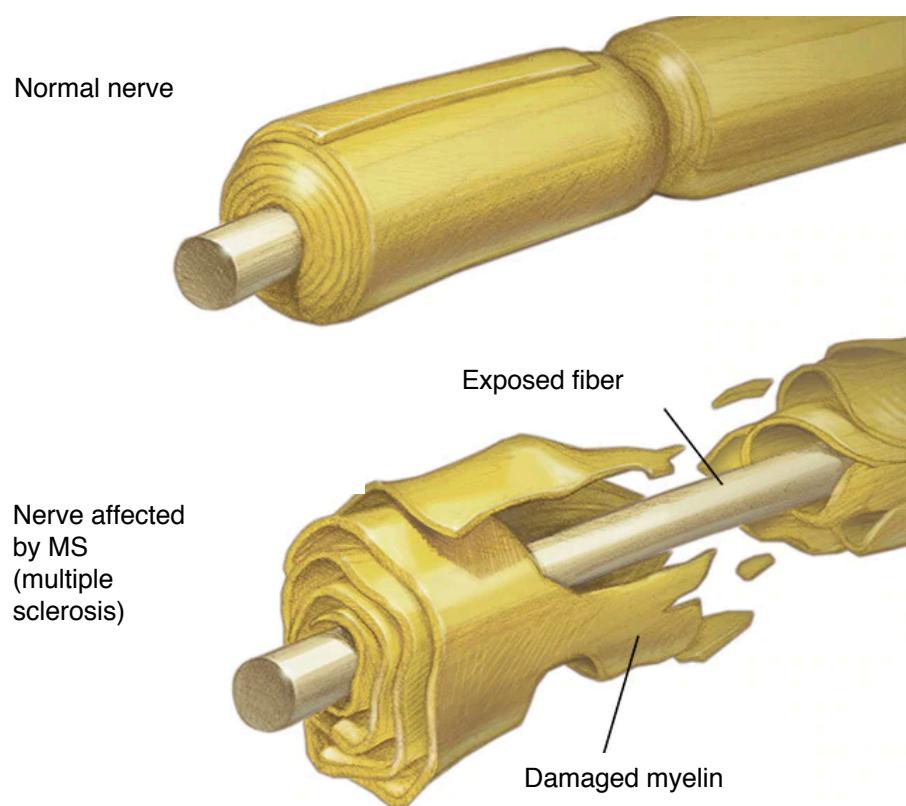
(Lin, Lei, Roeder: JASA, In revision)

Focus of this talk

Punch line (Statistics): eSVD performs a non-linear embedding by estimating a low-rank matrix of natural parameters via alternating minimization based on a random dot product model.

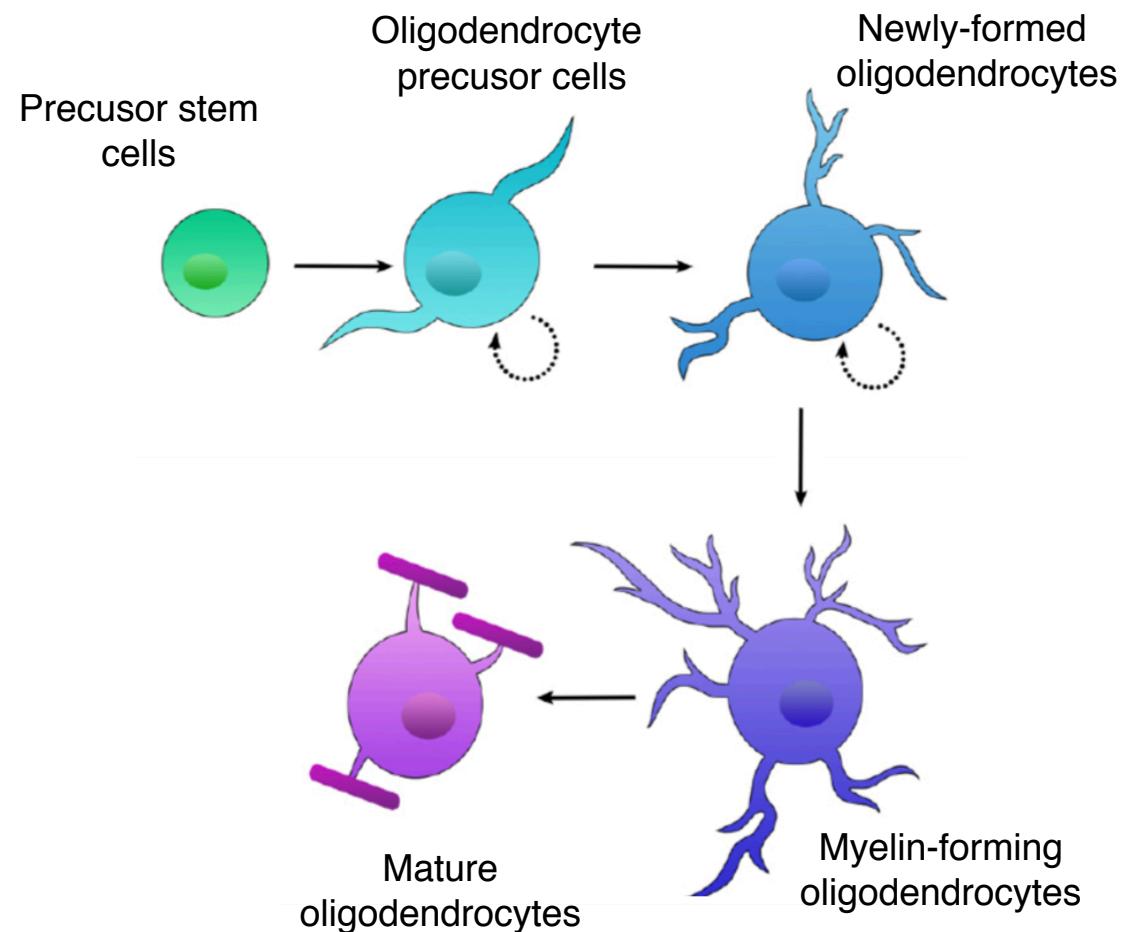
Oligodendrocytes are cells in the nervous system that produce myelin, and understanding their development can lead to better treatment for diseases such as MS.

- Oligodendrocytes
- Developmental trajectory
- Single-cell RNA-seq data
- Analysis pipeline



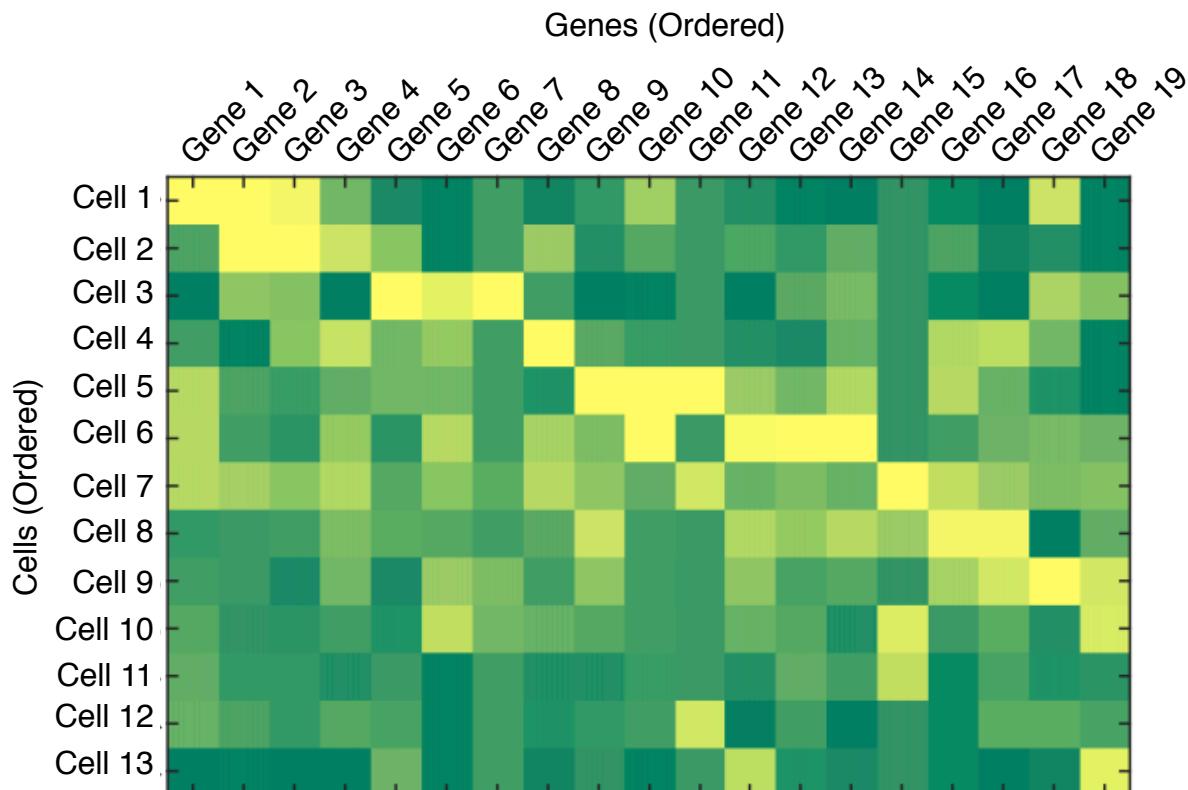
Oligodendrocytes are cells in the nervous system that produce myelin, and understanding their development can lead to better treatment for diseases such as MS.

- Oligodendrocytes
- Developmental trajectory
- Single-cell RNA-seq data
- Analysis pipeline



Oligodendrocytes are cells in the nervous system that produce myelin, and understanding their development can lead to better treatment for diseases such as MS.

- Oligodendrocytes
- Developmental trajectory
- Single-cell RNA-seq data
- Analysis pipeline



Oligodendrocytes are cells in the nervous system that produce myelin, and understanding their development can lead to better treatment for diseases such as MS.

Preprocessing → Embedding → Known cell types + Trajectory estimation

- Oligodendrocytes
- Developmental trajectory
- Single-cell RNA-seq data
- Analysis pipeline

Oligodendrocytes are cells in the nervous system that produce myelin, and understanding their development can lead to better treatment for diseases such as MS.

- Oligodendrocytes
- Developmental trajectory
- Single-cell RNA-seq data
- Analysis pipeline

Preprocessing → Embedding → Known cell types + Trajectory estimation

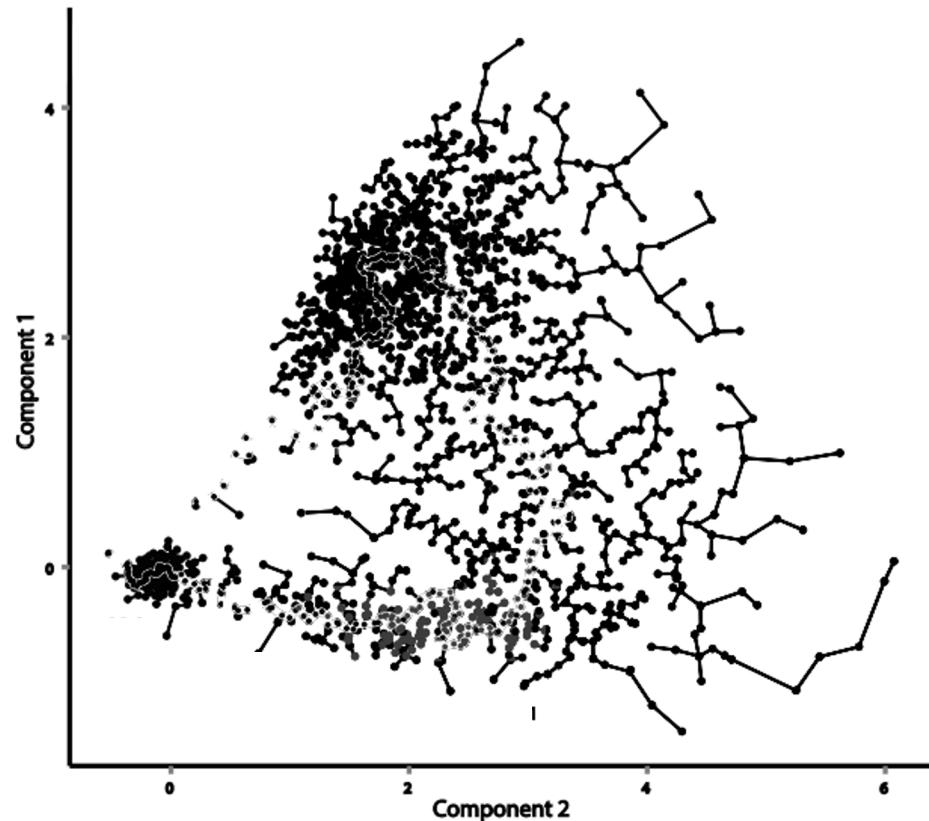


Figure from Marques et. al (2016)

Oligodendrocytes are cells in the nervous system that produce myelin, and understanding their development can lead to better treatment for diseases such as MS.

- Oligodendrocytes
- Developmental trajectory
- Single-cell RNA-seq data
- Analysis pipeline

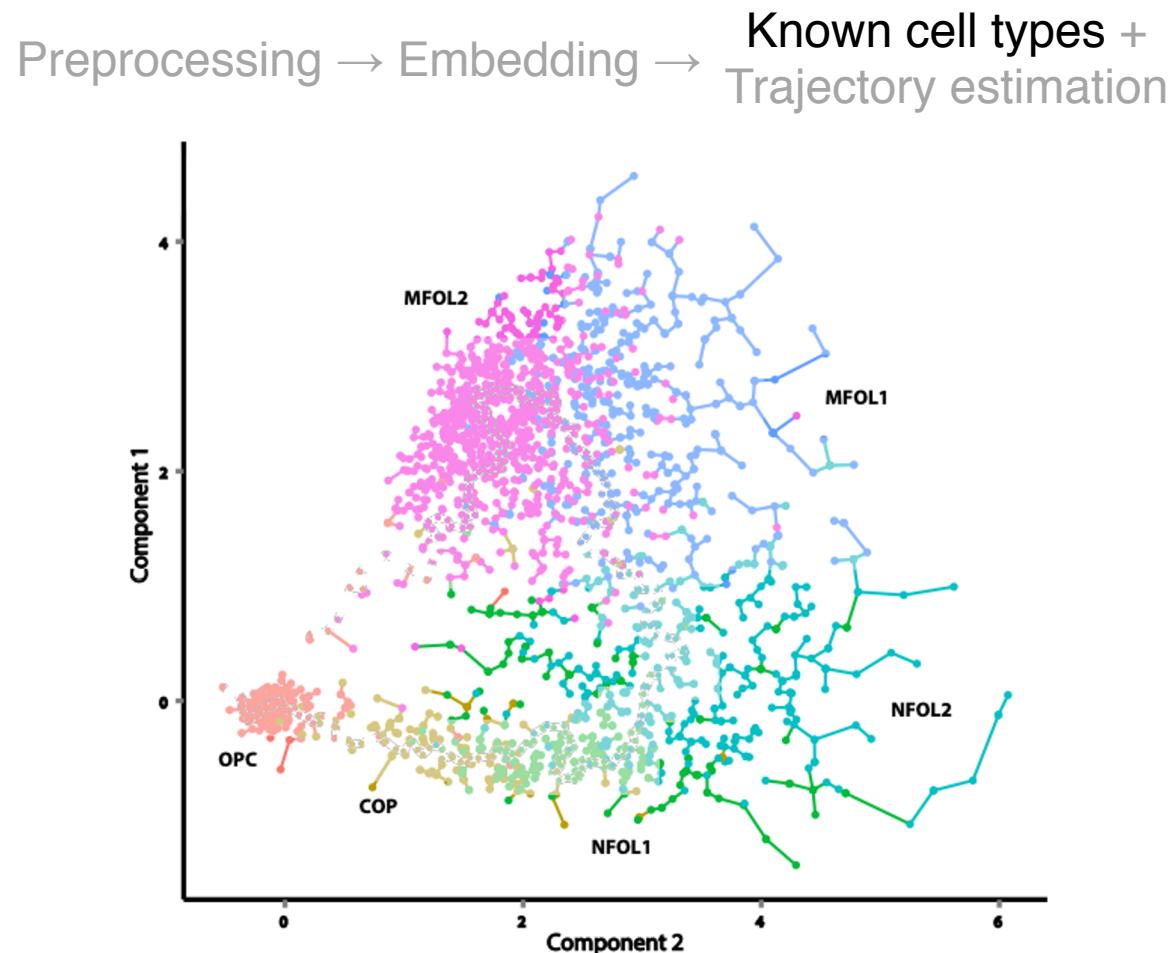


Figure from Marques et. al (2016)

Oligodendrocytes are cells in the nervous system that produce myelin, and understanding their development can lead to better treatment for diseases such as MS.

- Oligodendrocytes
- Developmental trajectory
- Single-cell RNA-seq data
- Analysis pipeline

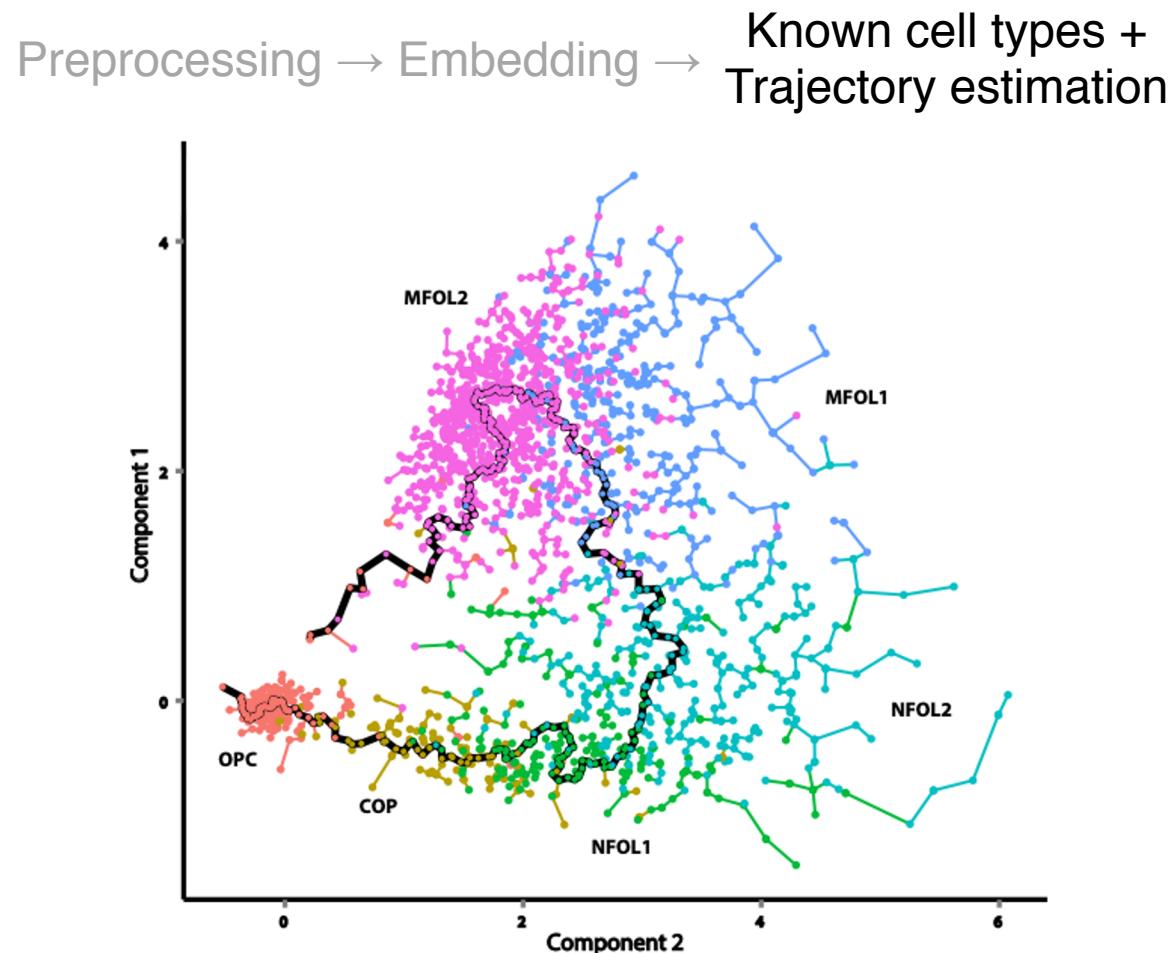


Figure from Marques et. al (2016)

Oligodendrocytes are cells in the nervous system that produce myelin, and understanding their development can lead to better treatment for diseases such as MS.

- Oligodendrocytes
- Developmental trajectory
- Single-cell RNA-seq data
- Analysis pipeline

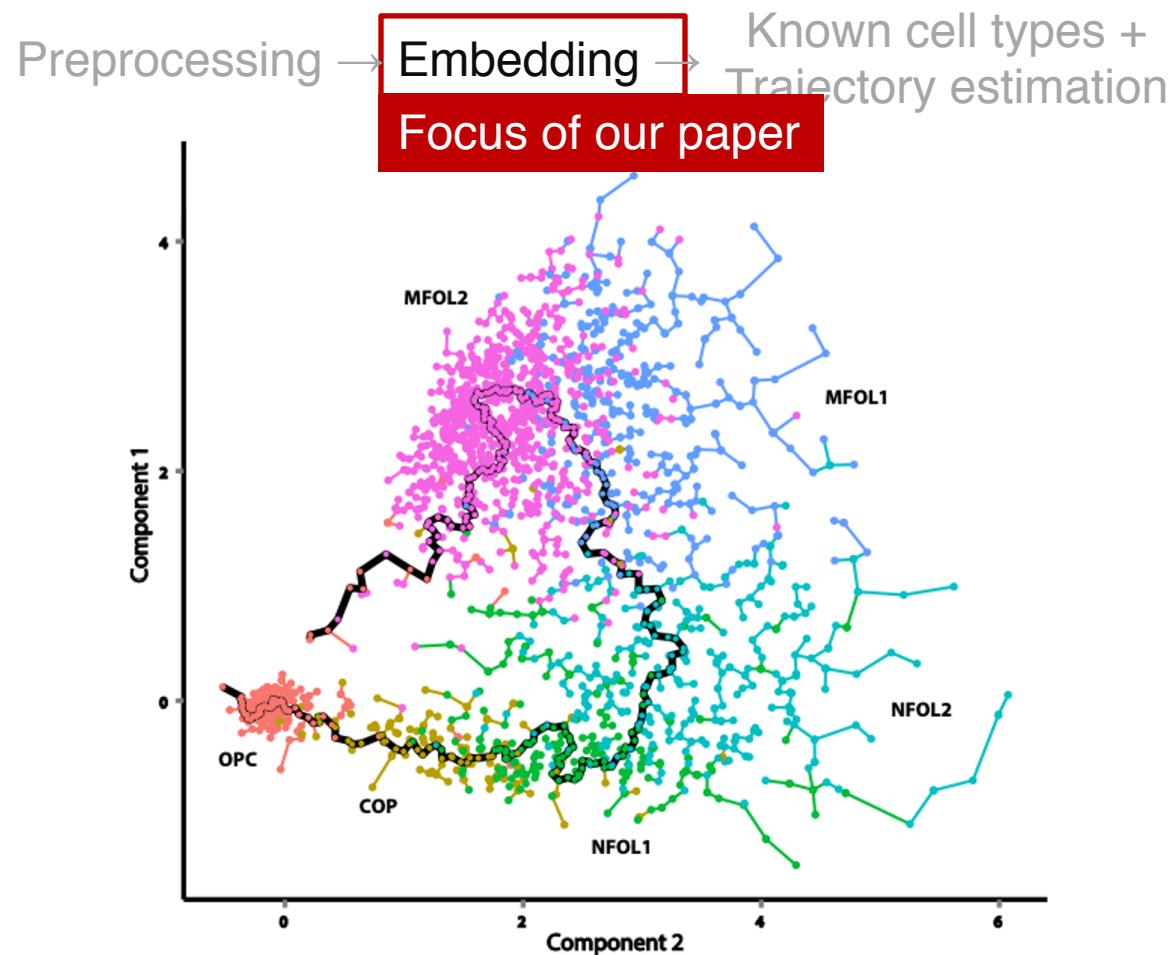


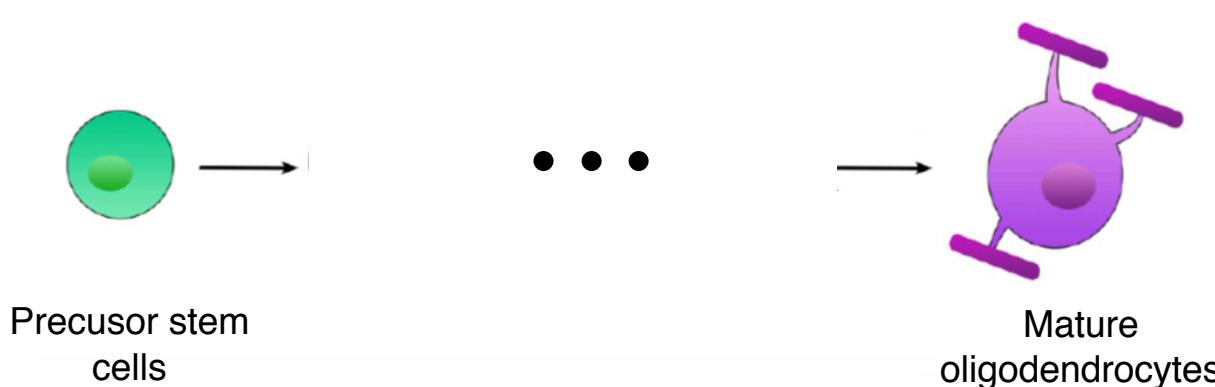
Figure from Marques et. al (2016)

The most common embedding, based on the SVD, is not desirable to model our single-cell RNA-seq dataset of oligodendrocytes from Marques et. al (2016).

$$A \in \mathbb{R}^{n \times p} : n \text{ cells and } p \text{ genes}$$

Preprocessed single-cell RNA-seq data

$n = 5000+$ cells
over 6 cell types,
 $p = 900+$ genes



The most common embedding, based on the SVD, is not desirable to model our single-cell RNA-seq dataset of oligodendrocytes from Marques et. al (2016).

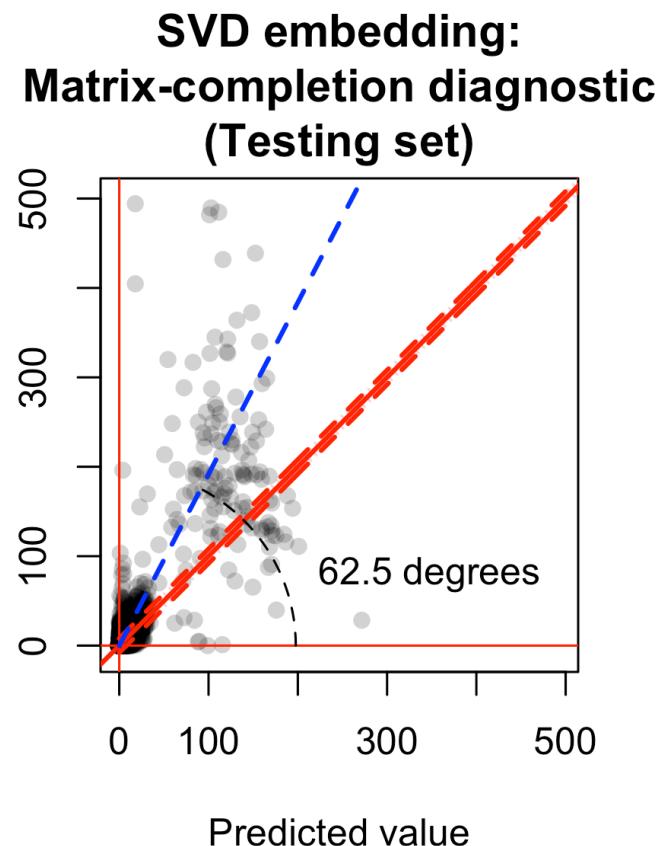
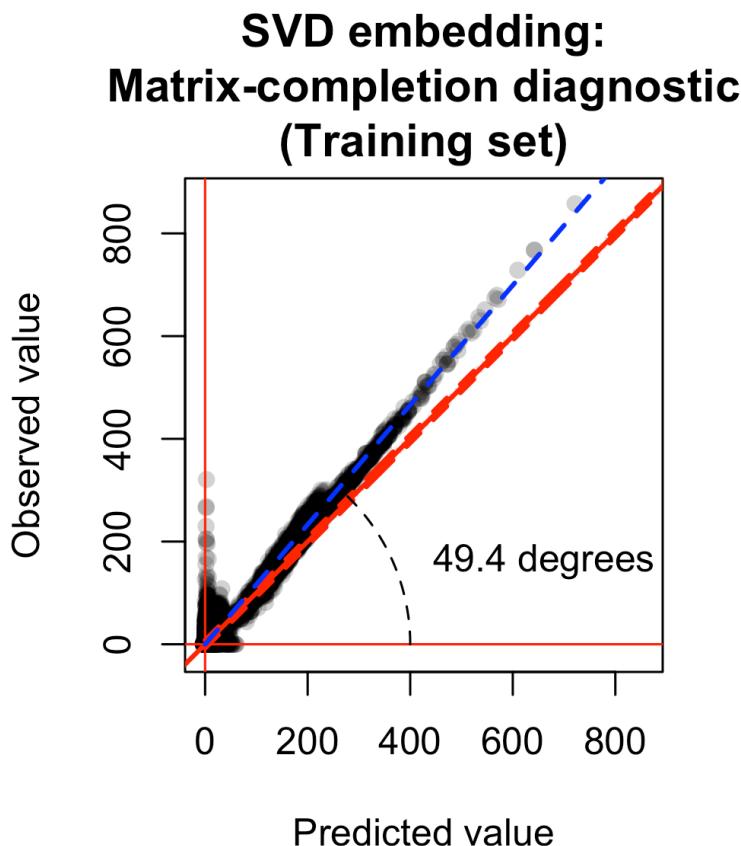
$A \in \mathbb{R}^{n \times p}$: n cells and p genes

$$A = \hat{U} \hat{D} \hat{V}^\top \quad \text{SVD}$$

$$\hat{X} = \left(\frac{n}{p}\right)^{1/4} \cdot \hat{U}_{1:k} \sqrt{\hat{D}_{1:k}}$$

$\mathbb{R}^{n \times k}$ matrix that encodes the lower-dimension embedding, based on a low-rank approximation of A

The most common embedding, based on the SVD, is not desirable to model our single-cell RNA-seq dataset of oligodendrocytes from Marques et. al (2016).



Using softImpute,
matrix-completion variant of SVD
(Mazumder et. al 2010)

The SVD embedding is actually a particular instance of a random dot product model, using the Gaussian distribution with constant variance.

Low-dim. embedding for cells

Low-dim. embedding for genes

$$\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{i.i.d.}{\sim} G,$$

and

$$\mathbf{Y}_1, \dots, \mathbf{Y}_p \stackrel{i.i.d.}{\sim} H$$

$$A_{ij} \sim F(\theta_{ij} = \mathbf{X}_i^\top \mathbf{Y}_j)$$

The SVD embedding is actually a particular instance of a random dot product model, using the Gaussian distribution with constant variance.

Low-dim. embedding for cells

Low-dim. embedding for genes

$$\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{i.i.d.}{\sim} G,$$

and

$$\mathbf{Y}_1, \dots, \mathbf{Y}_p \stackrel{i.i.d.}{\sim} H$$

$$A_{ij} \sim F(\theta_{ij} = \mathbf{X}_i^\top \mathbf{Y}_j)$$

Preprocessed single-cell
RNA-seq data

A single-parameter
distribution family

The SVD embedding is actually a particular instance of a random dot product model, using the Gaussian distribution with constant variance.

Low-dim. embedding for cells

Low-dim. embedding for genes

$$\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{i.i.d.}{\sim} G,$$

and

$$\mathbf{Y}_1, \dots, \mathbf{Y}_p \stackrel{i.i.d.}{\sim} H$$

$$A_{ij} \sim F(\theta_{ij} = \underbrace{\mathbf{X}_i^\top \mathbf{Y}_j}_{})$$

Preprocessed single-cell
RNA-seq data

Dot product between
latent vectors

The SVD embedding is actually a particular instance of a random dot product model, using the Gaussian distribution with constant variance.

Low-dim. embedding for cells

$$\mathbf{X}_1, \dots, \mathbf{X}_n \stackrel{i.i.d.}{\sim} G, \quad \text{and} \quad \mathbf{Y}_1, \dots, \mathbf{Y}_p \stackrel{i.i.d.}{\sim} H$$

$$A_{ij} \sim F(\theta_{ij} = \mathbf{X}_i^\top \mathbf{Y}_j)$$

Exponential-family distribution

$$p(A_{ij} \mid \theta_{ij}) \propto \exp(A_{ij} \cdot \theta_{ij} - g(\theta_{ij}))$$

Natural parameter

The SVD embedding is actually a particular instance of a random dot product model, using the Gaussian distribution with constant variance.

Low-dim. embedding for cells

$$\{X_1, \dots, X_n\} \stackrel{i.i.d.}{\sim} G, \quad \text{and} \quad \{Y_1, \dots, Y_p\} \stackrel{i.i.d.}{\sim} H$$

$$A_{ij} \sim F(\theta_{ij} = X_i^\top Y_j)$$

Exponential-family distribution

Intuitive idea: For a given one-parameter exponential-family distribution F , estimate the embedding $\{X_1, \dots, X_n\}$ by minimize the negative log-likelihood using A (i.e., “MLE”)

The SVD embedding is actually a particular instance of a random dot product model, using the Gaussian distribution with constant variance.

Low-dim. embedding for cells

$$\{X_1, \dots, X_n\} \stackrel{i.i.d.}{\sim} G, \quad \text{and} \quad \{Y_1, \dots, Y_p\} \stackrel{i.i.d.}{\sim} H$$

$$A_{ij} \sim F(\theta_{ij} = X_i^\top Y_j)$$

$$\mathcal{L}(X, Y) = \frac{1}{np} \sum_{(i,j)} \left[g(X_i^\top Y_j) - A_{ij} \cdot (X_i^\top Y_j) \right]$$

Loss function
to minimize

Negative log-likelihood

The SVD embedding is actually a particular instance of a random dot product model, using the Gaussian distribution with constant variance.

Low-dim. embedding for cells

$$\{X_1, \dots, X_n\} \stackrel{i.i.d.}{\sim} G, \quad \text{and} \quad \{Y_1, \dots, Y_p\} \stackrel{i.i.d.}{\sim} H$$

$$A_{ij} \sim F(\theta_{ij} = X_i^\top Y_j)$$

$$\mathcal{L}(X, Y) = \frac{1}{np} \sum_{(i,j)} \left[g(X_i^\top Y_j) - A_{ij} \cdot (X_i^\top Y_j) \right]$$

Plugging in
Gaussian w/
constant
variance

$$\propto \frac{1}{np} \sum_{(i,j)} \left[\left(A_{ij} - (X_i^\top Y_j) \right)^2 \right]$$

(Solved via SVD)

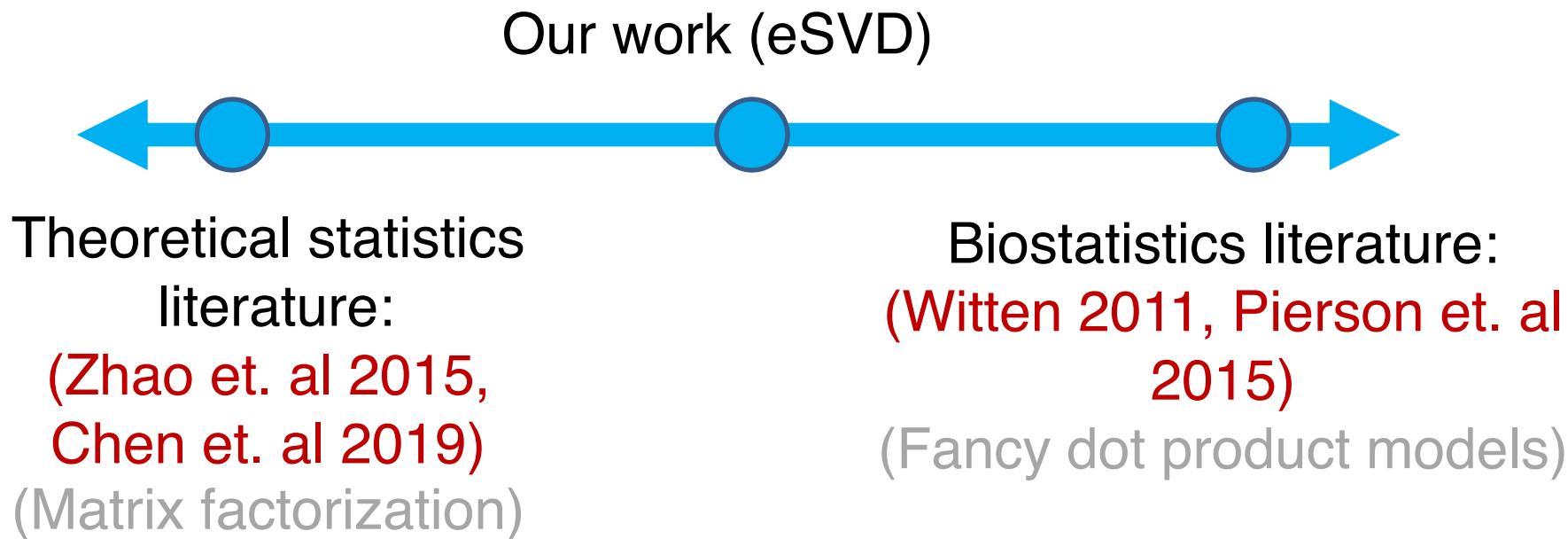
Many previous papers that use similar formulations either lack modeling flexibility or the desired theoretical properties.



Theoretical statistics
literature:
(Zhao et. al 2015,
Chen et. al 2019)
(Matrix factorization)

Biostatistics literature:
(Witten 2011, Pierson et. al
2015)
(Fancy dot product models)

Many previous papers that use similar formulations either lack modeling flexibility or the desired theoretical properties.



Goals: Computational efficiency, theoretical foundation, modeling flexibility, refinement of random dot product model, amendable for a tuning procedure

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

$$\arg \min_{\mathbf{X}, \mathbf{Y}} \mathcal{L}(\mathbf{X}, \mathbf{Y}) = \arg \min_{\mathbf{X}, \mathbf{Y}} \frac{1}{np} \sum_{(i,j)} \left[g(\mathbf{X}_i^\top \mathbf{Y}_j) - A_{ij} \cdot (\mathbf{X}_i^\top \mathbf{Y}_j) \right]$$

$$\mathbf{X} \in \mathbb{R}^{n \times k}, \mathbf{Y} \in \mathbb{R}^{p \times k} \quad (\text{where } \mathbf{XY}^\top = \Theta)$$

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

$$\arg \min_{\mathbf{X}, \mathbf{Y}} \mathcal{L}(\mathbf{X}, \mathbf{Y}) = \arg \min_{\mathbf{X}, \mathbf{Y}} \frac{1}{np} \sum_{(i,j)} \left[g(\mathbf{X}_i^\top \mathbf{Y}_j) - A_{ij} \cdot (\mathbf{X}_i^\top \mathbf{Y}_j) \right]$$

- Stage 1: Initialization
 - Initialize the estimate based on SVD of some transformation of A

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

$$\arg \min_{\mathbf{X}, \mathbf{Y}} \mathcal{L}(\mathbf{X}, \mathbf{Y}) = \arg \min_{\mathbf{X}, \mathbf{Y}} \frac{1}{np} \sum_{(i,j)} \left[g(\mathbf{X}_i^\top \mathbf{Y}_j) - A_{ij} \cdot (\mathbf{X}_i^\top \mathbf{Y}_j) \right]$$

- Stage 1: Initialization
 - Initialize the estimate based on SVD of some transformation of A
- Stage 2: Alternating minimization
 - Update the estimate by re-estimating X , keeping Y fixed, or re-estimating Y , keeping X fixed

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

$$\arg \min_{\mathbf{X}, \mathbf{Y}} \mathcal{L}(\mathbf{X}, \mathbf{Y}) = \arg \min_{\mathbf{X}, \mathbf{Y}} \frac{1}{np} \sum_{(i,j)} \left[g(\mathbf{X}_i^\top \mathbf{Y}_j) - A_{ij} \cdot (\mathbf{X}_i^\top \mathbf{Y}_j) \right]$$

- Stage 1: Initialization
 - Initialize the estimate based on SVD of some transformation of A
- Stage 2: Alternating minimization
 - Update the estimate by re-estimating \mathbf{X} , keeping \mathbf{Y} fixed, or re-estimating \mathbf{Y} , keeping \mathbf{X} fixed
- Stage 3: Reparameterization
 - Ensure identifiability conditions are met; $\left(\frac{1}{n} \hat{\mathbf{X}}^\top \hat{\mathbf{X}} = \frac{1}{p} \hat{\mathbf{Y}}^\top \hat{\mathbf{Y}} \right)$

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

- Stage 2: Alternating minimization

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

- Stage 2: Alternating minimization

$$\mathbf{X}^{(t)} = \underset{\mathbf{X} \in \mathbb{R}^{n \times k}}{\operatorname{argmin}} \mathcal{L}(\mathbf{X}, \bar{\mathbf{Y}}^{(t)}) : \mathbf{X}_i^\top \bar{\mathbf{Y}}_j^{(t)} \in \mathcal{R}, \quad \forall(i, j),$$

$$\bar{\mathbf{X}}^{(t+1)} = \sqrt{n} \cdot \text{LeftSVD}(\mathbf{X}^{(t)}),$$

$$\mathbf{Y}^{(t+1)} = \underset{\mathbf{Y} \in \mathbb{R}^{p \times k}}{\operatorname{argmin}} \mathcal{L}(\bar{\mathbf{X}}^{(t+1)}, \mathbf{Y}) : (\bar{\mathbf{X}}_i^{(t+1)})^\top \mathbf{Y}_j \in \mathcal{R}, \quad \forall(i, j),$$

$$\bar{\mathbf{Y}}^{(t+1)} = \sqrt{p} \cdot \text{LeftSVD}(\mathbf{Y}^{(t+1)})$$

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

- Stage 2: Alternating minimization

$$\mathbf{X}^{(t)} = \underset{\mathbf{X} \in \mathbb{R}^{n \times k}}{\operatorname{argmin}} \mathcal{L}(\mathbf{X}, \bar{\mathbf{Y}}^{(t)})$$

$$\mathbf{X}^{(t+1)} = \sqrt{n} \cdot \text{LeftSVD}(\mathbf{X}^{(t)}),$$

$$\mathbf{Y}^{(t+1)} = \underset{\mathbf{Y} \in \mathbb{R}^{p \times k}}{\operatorname{argmin}} \mathcal{L}(\bar{\mathbf{X}}^{(t+1)}, \mathbf{Y})$$

$$\bar{\mathbf{Y}}^{(t+1)} = \sqrt{p} \cdot \text{LeftSVD}(\mathbf{Y}^{(t+1)})$$

Minimizing one of the factors based on the log-likelihood (Convex problem)

$\forall (i, j),$

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

- Stage 2: Alternating minimization

Need to ensure the inner product still yields a valid natural parameter

We use Frank-Wolfe to do this constrained optimization

$$: \quad \mathbf{X}_i^\top \bar{\mathbf{Y}}_j^{(t)} \in \mathcal{R}, \quad \forall(i, j),$$

$$), \\ \curvearrowright : (\bar{\mathbf{X}}_i^{(t+1)})^\top \mathbf{Y}_j \in \mathcal{R}, \quad \forall(i, j), \\ \cdots \\ \curvearrowleft : (\bar{\mathbf{Y}}_j^{(t+1)})^\top \mathbf{X}_i \in \mathcal{R}, \quad \forall(i, j),$$

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

- Stage 2: Alternating minimization

$$\mathbf{X}^{(t)} = \underset{\mathbf{X} \in \mathbb{R}^{n \times k}}{\operatorname{argmin}} \mathcal{L}(\mathbf{X}, \bar{\mathbf{Y}}^{(t)}) : \mathbf{X}_i^\top \bar{\mathbf{Y}}_j^{(t)} \in \mathcal{R}, \quad \forall (i, j),$$

$$\bar{\mathbf{X}}^{(t+1)} = \sqrt{n} \cdot \text{LeftSVD}(\mathbf{X}^{(t)}),$$

$$\mathbf{Y}^{(t+1)} = \underset{\mathbf{Y} \in \mathbb{R}^{p \times k}}{\operatorname{argmin}} \mathcal{L}(\bar{\mathbf{X}}^{(t+1)}, \mathbf{Y})$$

$$\bar{\mathbf{Y}}^{(t+1)} = \sqrt{p} \cdot \text{LeftSVD}(\mathbf{Y}^{(t+1)})$$

We reparameterize each estimate by its left singular vectors for identifiability reasons

eSVD performs a suitable initialization then uses alternating minimization based on the negative log-likelihood.

- Stage 2: Alternating minimization

$$\mathbf{X}^{(t)} = \underset{\mathbf{X} \in \mathbb{R}^{n \times k}}{\operatorname{argmin}} \mathcal{L}(\mathbf{X}, \bar{\mathbf{Y}}^{(t)}) : \mathbf{X}_i^\top \bar{\mathbf{Y}}_j^{(t)} \in \mathcal{R}, \quad \forall(i, j),$$

$$\bar{\mathbf{X}}^{(t+1)} = \sqrt{n} \cdot \text{LeftSVD}(\mathbf{X}^{(t)}),$$

$$\mathbf{Y}^{(t+1)} = \underset{\mathbf{Y} \in \mathbb{R}^{p \times k}}{\operatorname{argmin}} \mathcal{L}(\bar{\mathbf{X}}^{(t+1)}, \mathbf{Y}) :$$

$$\bar{\mathbf{Y}}^{(t+1)} = \sqrt{p} \cdot \text{LeftSVD}(\mathbf{Y}^{(t+1)})$$

We then move to the next factor. We keep alternating between the two factors until convergence.

We prove convergence rates using theoretical techniques from nonconvex optimization and network literature.

- Estimating low-rank matrix of natural parameters:
 - Requires understanding the sufficient conditions for initialization as well as considering the minimization on a population level (Zhao et al. 2015, Balakrishnan et al. 2017)

We prove convergence rates using theoretical techniques from nonconvex optimization and network literature.

- Estimating low-rank matrix of natural parameters:
 - Requires understanding the sufficient conditions for initialization as well as considering the minimization on a population level ([Zhao et al. 2015](#), [Balakrishnan et al. 2017](#))
- Handling the extra source of randomness due to the hierarchical model:
 - Requires careful study of identifiability conditions and perturbation bounds ([Lei 2019](#))

We prove convergence rates using theoretical techniques from nonconvex optimization and network literature.

- Estimating low-rank matrix of natural parameters:
 - Requires understanding the sufficient conditions for initialization as well as considering the minimization on a population level ([Zhao et al. 2015](#), [Balakrishnan et al. 2017](#))
- Handling the extra source of randomness due to the hierarchical model:
 - Requires careful study of identifiability conditions and perturbation bounds ([Lei 2019](#))
- Applying these generic bounds to specific distribution families:
 - Requires deploying empirical process and concentration bounds ([Candes and Sur 2019](#))

eSVD with a curved Gaussian distribution models the oligodendrocyte dataset much more appropriately.

$$A_{ij} \sim F = \text{Normal}(\mu_{ij}, \mu_{ij}^2/4)$$

so the natural parameter becomes:

$$\mathbf{X}_i^\top \mathbf{Y}_j = \theta_{ij} = -1/\mu_{ij}$$

- Most of the distribution's mass is positive
- Overall, allows for more variability compared to the negative binomial or Poisson distribution

eSVD with a curved Gaussian distribution models the oligodendrocyte dataset much more appropriately.

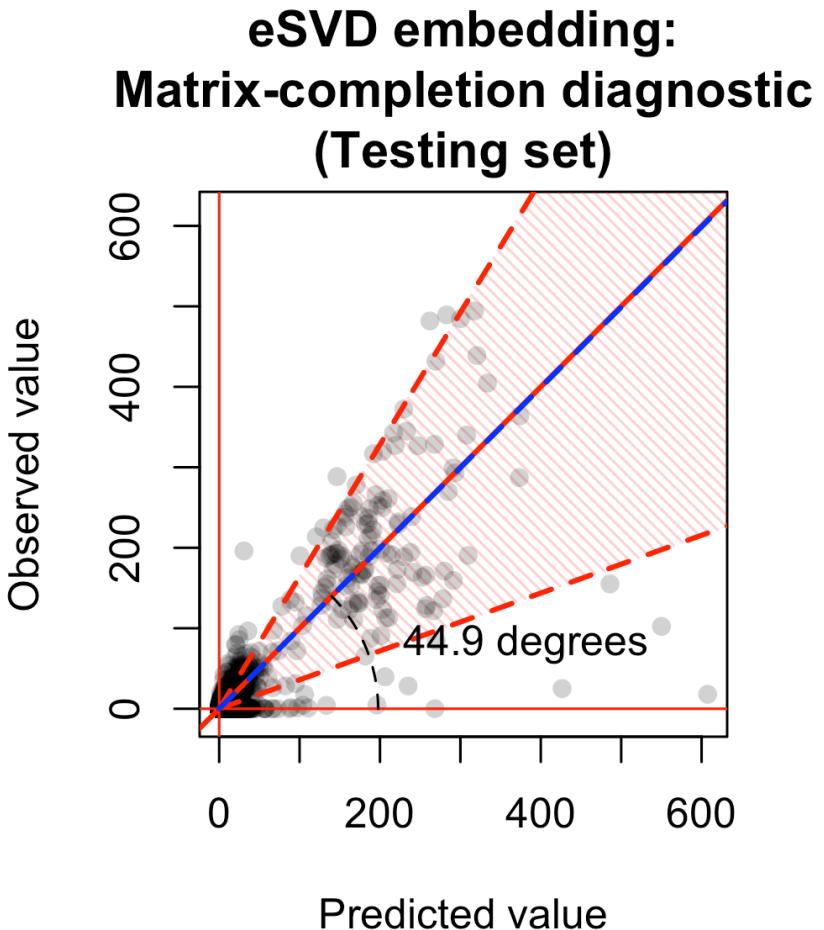
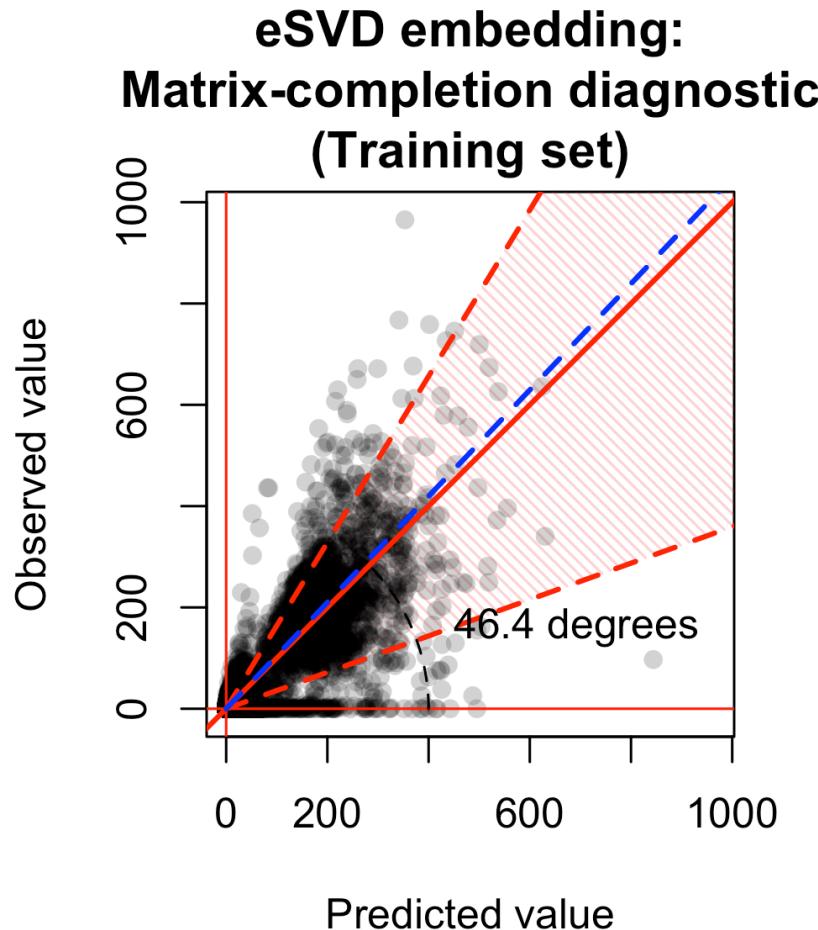
$$A_{ij} \sim F = \text{Normal}(\mu_{ij}, \mu_{ij}^2/4)$$

so the natural parameter becomes:

$$\mathbf{X}_i^\top \mathbf{Y}_j = \theta_{ij} = -1/\mu_{ij}$$

- Most of the distribution's mass is positive
- Overall, allows for more variability compared to the negative binomial or Poisson distribution
- The constant 4 or the dimensionality of the latent space k can be tuned via the matrix-completion diagnostic

eSVD with a curved Gaussian distribution models the oligodendrocyte dataset much more appropriately.



The eSVD embedding enables us to discover two branching trajectories in downstream analyses, whereas the SVD embedding only leads to one trajectory.

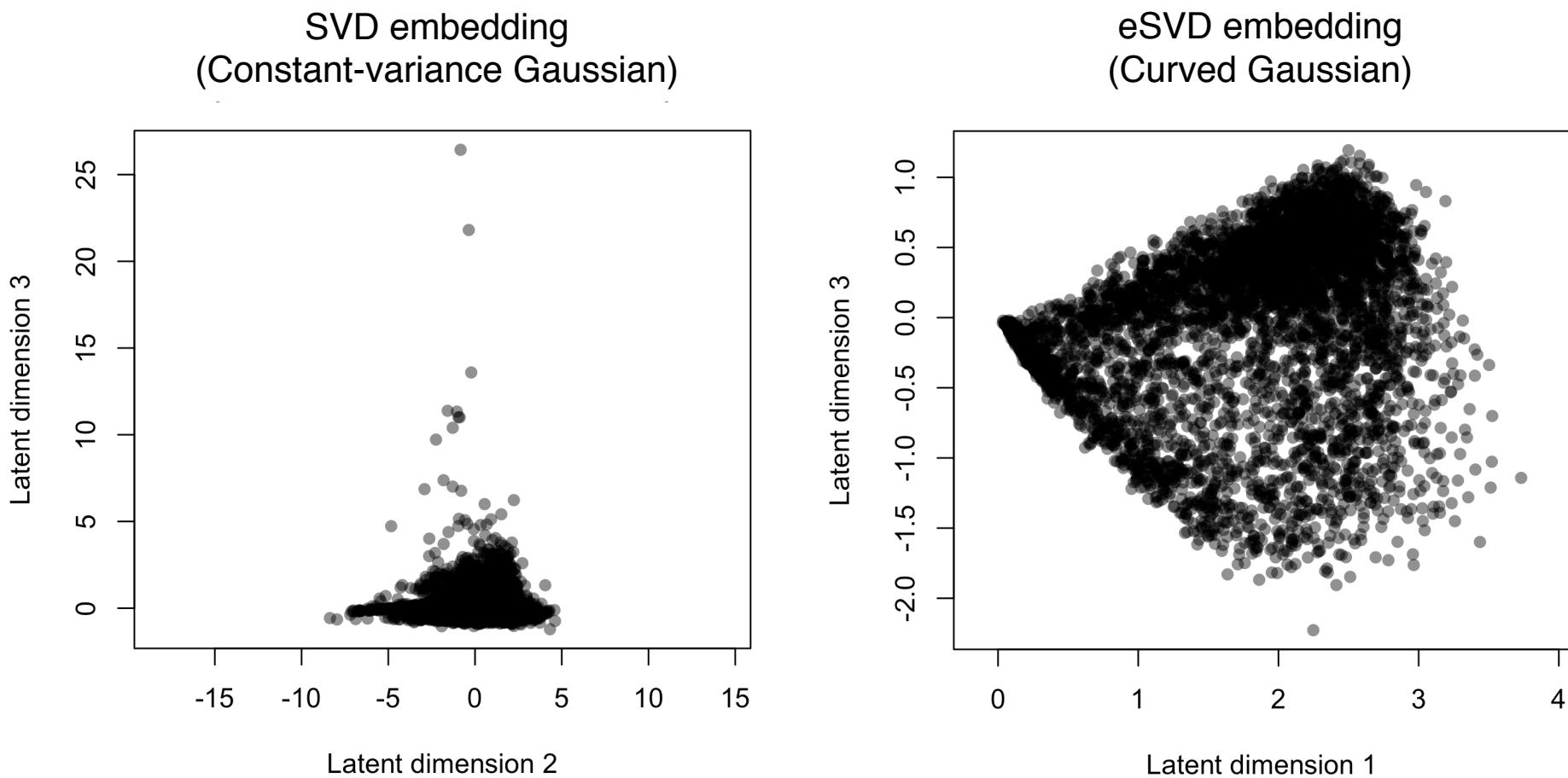
Color coding
in upcoming
plots

Major Cell Type (6 total)	Number of cell sub-types (13 total)
Pdgfra+ precursor	1
Oligodendrocyte precursor cell	1
Differentiation-committed oligodendrocyte precursors	1
Newly formed oligodendrocytes	2
Myelin-forming oligodendrocytes	2
Mature oligodendrocyte	6



Reminder: $n = 5000+$ cells in total and $p = 900+$ genes

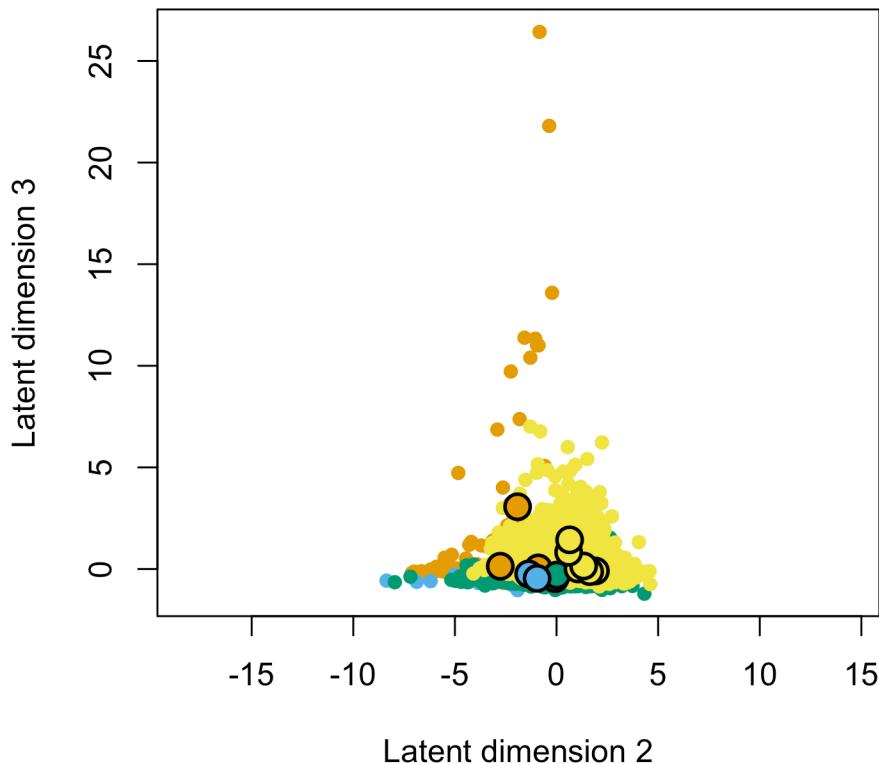
The eSVD embedding enables us to discover two branching trajectories in downstream analyses, whereas the SVD embedding only leads to one trajectory.



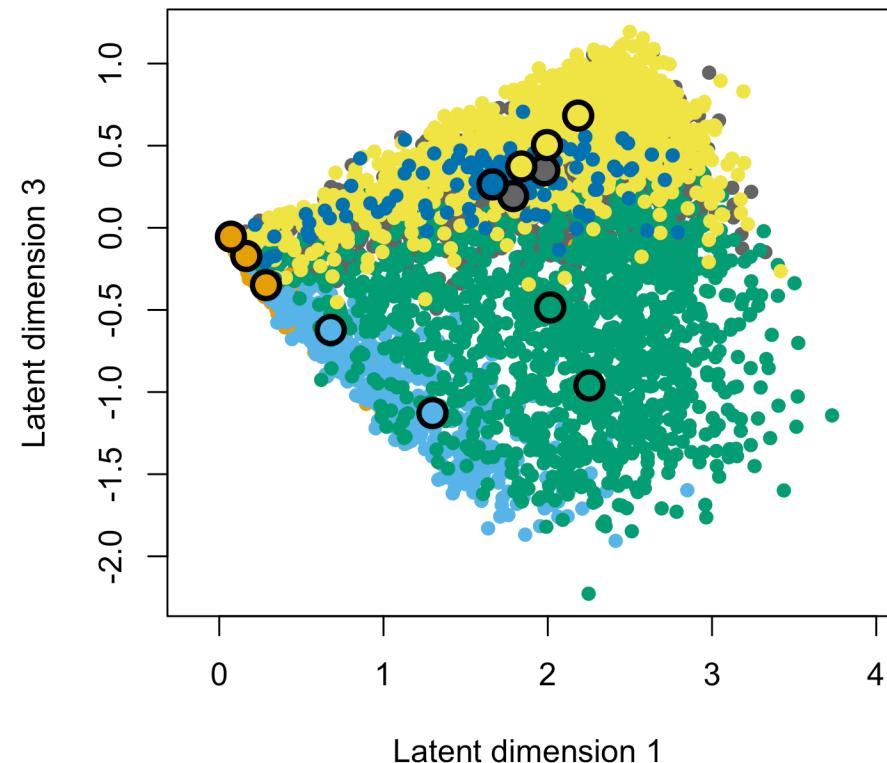
Preprocessing → **Embedding** → Known cell types +
Trajectory estimation

The eSVD embedding enables us to discover two branching trajectories in downstream analyses, whereas the SVD embedding only leads to one trajectory.

SVD embedding
(Constant-variance Gaussian)

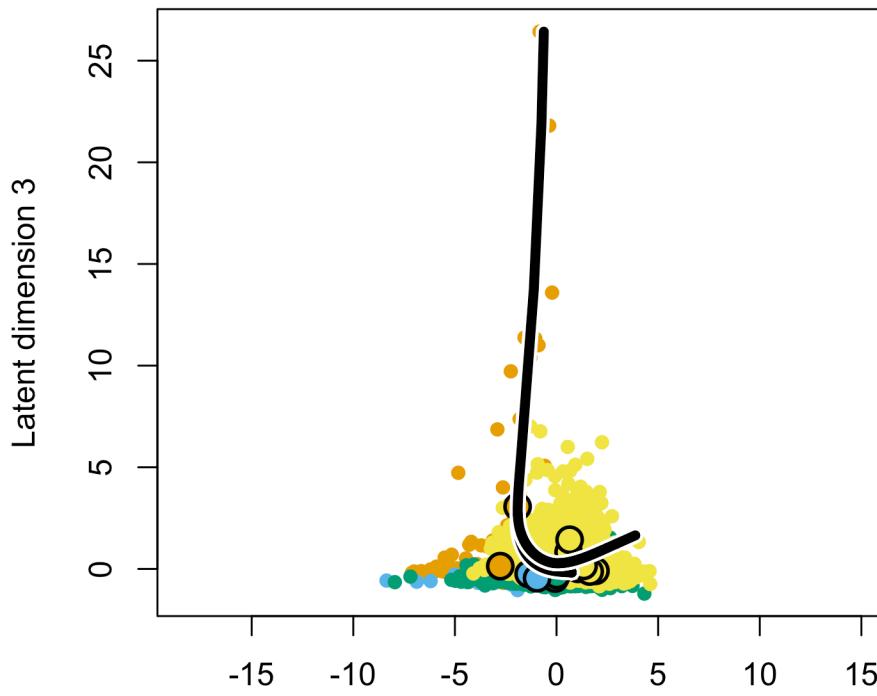


eSVD embedding
(Curved Gaussian)



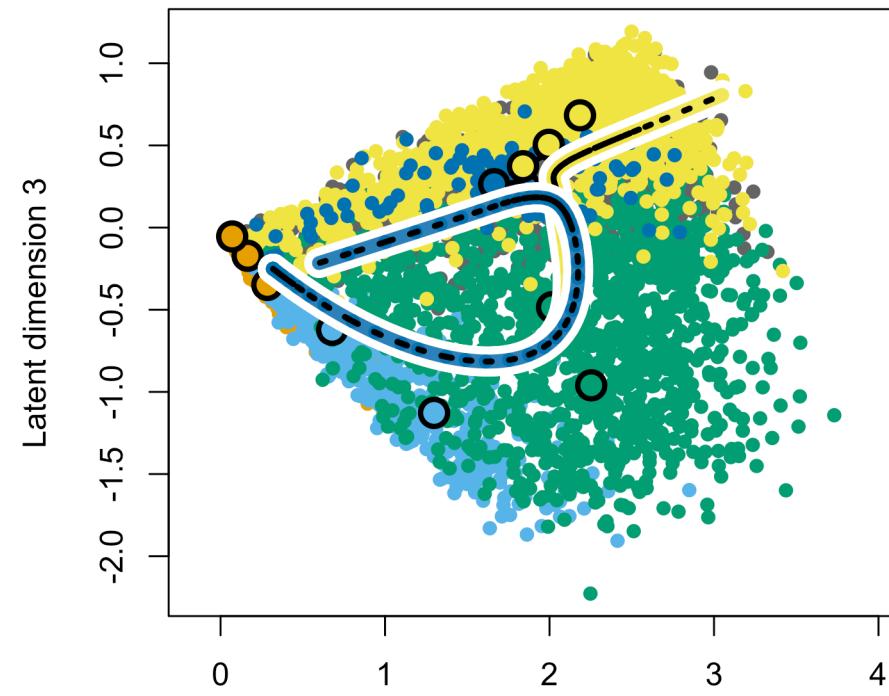
The eSVD embedding enables us to discover two branching trajectories in downstream analyses, whereas the SVD embedding only leads to one trajectory.

SVD embedding and trajectory
(Constant-variance Gaussian)



Using Slingshot,
a trajectory estimator
(Street et. al 2018)

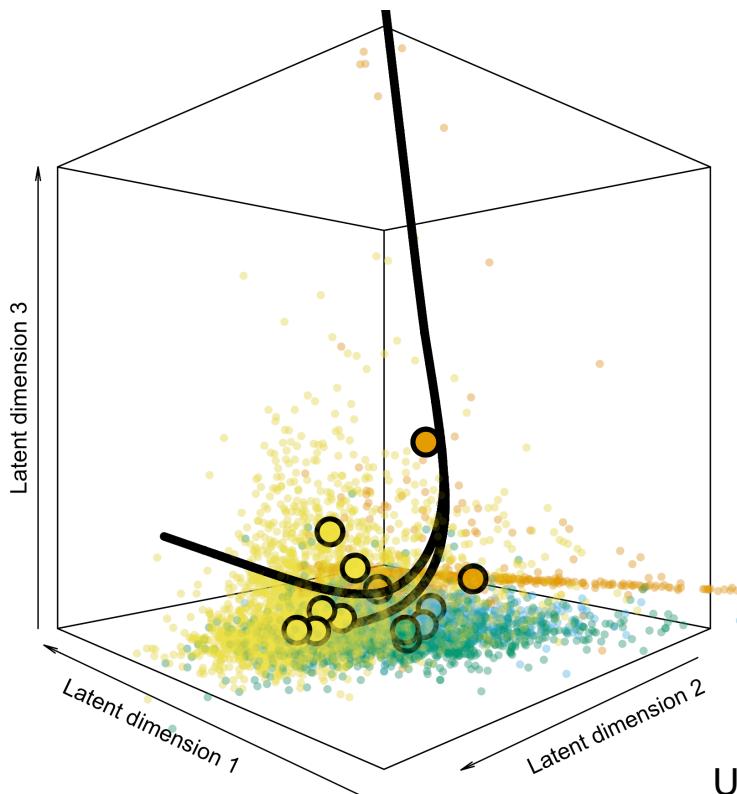
eSVD embedding and trajectory
(Curved Gaussian)



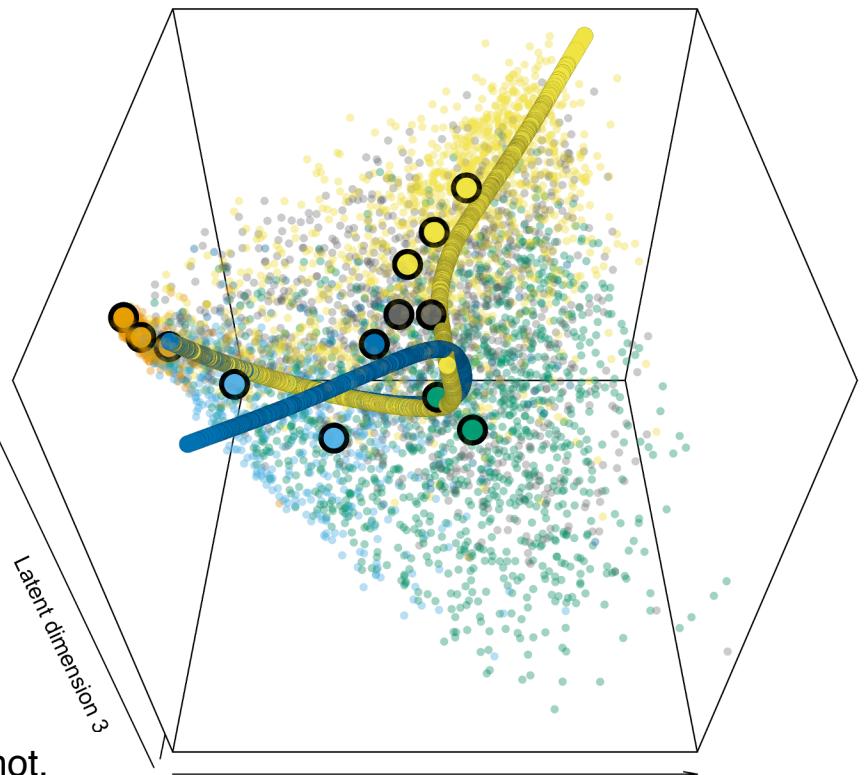
Latent dimension 1

The eSVD embedding enables us to discover two branching trajectories in downstream analyses, whereas the SVD embedding only leads to one trajectory.

SVD embedding and trajectory
(Constant-variance Gaussian)



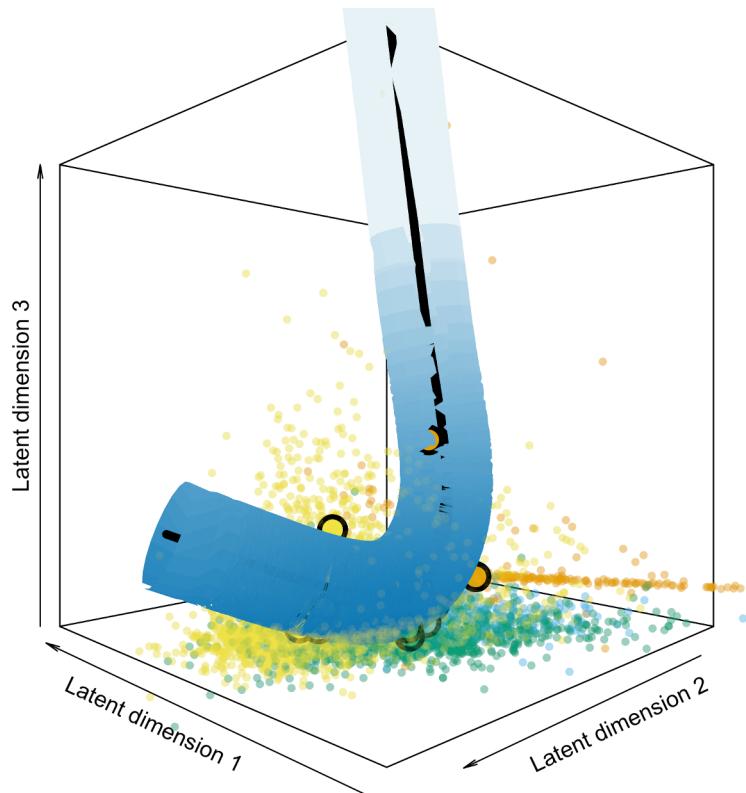
eSVD embedding and trajectory
(Curved Gaussian)



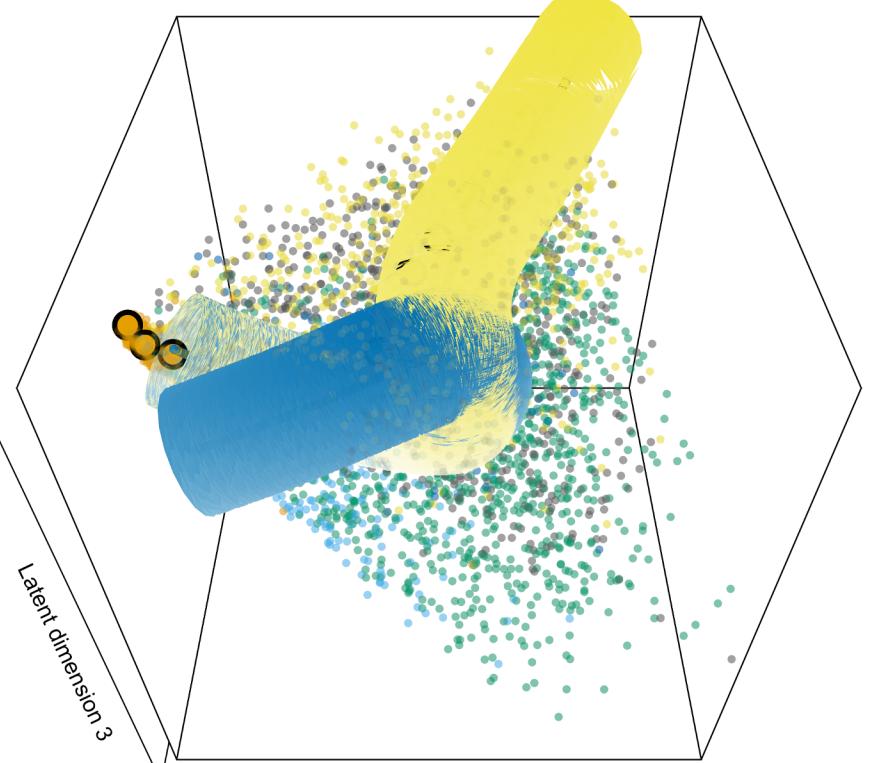
Using Slingshot,
a trajectory estimator
(Street et. al 2018)

The eSVD embedding enables us to discover two branching trajectories in downstream analyses, whereas the SVD embedding only leads to one trajectory.

SVD embedding and uncertainty tube
(Constant-variance Gaussian)



eSVD embedding and uncertainty tube
(Curved Gaussian)



Multiple trajectories in agreement with ongoing research
(Marques et. al 2018, Jakel et. al 2019)

Since embeddings offer many more downstream applications, future work can strengthen eSVD's theoretical properties as well as modeling flexibility.

Other potential downstream applications:

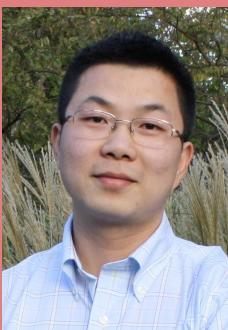
- Visualization
- Imputation
- Clustering (or bi-clustering)

Summary: eSVD is a non-linear embedding method based on any one-parameter exponential-family distribution, and we prove theoretical foundation and apply it to study oligodendrocytes' developmental trajectories.

Thank you!



Kathryn
Roeder



Jing
Lei



Max
G'Sell



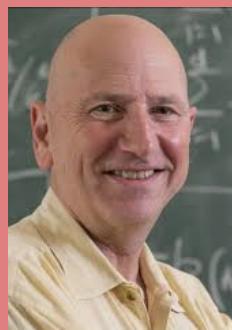
Han
Liu



Anjali
Mazumder



Ryan J.
Tibshirani



Larry
Wasserman

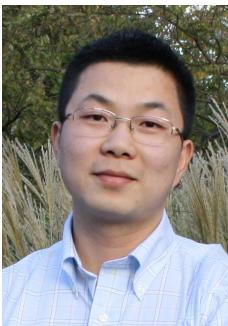
eSVD is a non-linear embedding method based on any one-parameter exponential-families, and we prove theoretical foundation and apply it to study oligodendrocytes' developmental trajectories.

Possible questions to ask:

- More details about oligodendrocyte data or known scientific findings
- Identifiability condition
- Assumptions for convergence
- Theorems for convergence
- Details on this missing-value diagnostic/tuning procedure
- Results from simulation studies
- Details about how Slingshot works (i.e., how to estimate developmental trajectories)
- Relation to other non-linear embedding methods or matrix factorization work
- Thoughts about how to improve eSVD in future work
- Brief overview of the other chapters in thesis



Kathryn
Roeder



Jing
Lei



Max
G'Sell



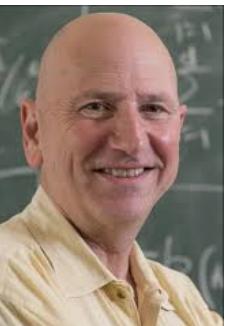
Han
Liu



Anjali
Mazumder



Ryan J.
Tibshirani



Larry
Wasserman