tvannoy.github.io/randomized-subkmeans-presentation

# IMPROVED SUBSPACE K-MEANS PERFORMANCE VIA A RANDOMIZED MATRIX DECOMPOSITION

**Trevor Vannoy**

Jacob Senecal, Veronika Strnadova-Neeley

# OUTLINE

- Clustering in high dimensions

- Subspace k-means

- Randomized Subspace k-means

- Experiments

- Results

- Conclusions

# Clustering in High Dimensions

# PROBLEMS

- curse of dimensionality
- computational complexity
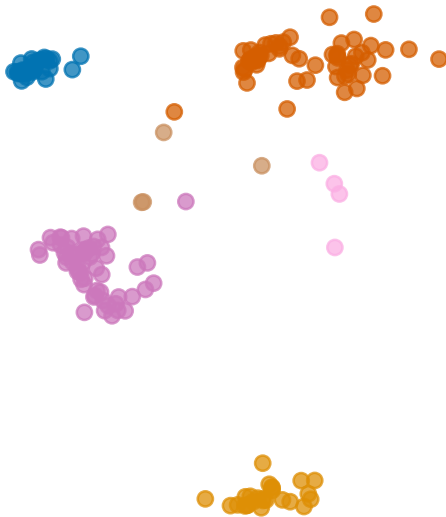- visualization

# SOLUTIONS

- subspace clustering
  - separate subspaces
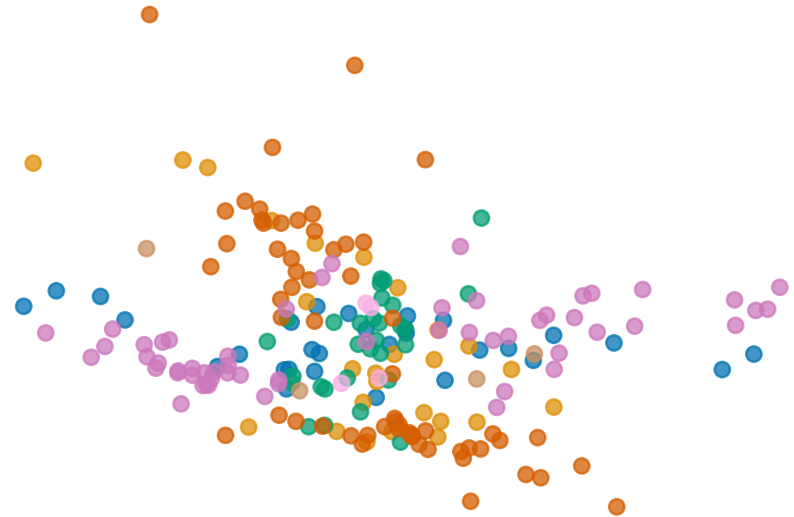  - common subspace

# Subspace k-means

# SUBSPACE K-MEANS

- transforms data into a cluster subspace and a noise subspace

- alternates between subspace estimation and clustering

**cluster subspace**　　　**noise subspace**

# OBJECTIVE FUNCTION

$$\mathcal{J} = \left[ \sum_{i=1}^{k} \sum_{\mathbf{x} \in C_i} \| \quad \mathbf{x} - \quad \boldsymbol{\mu}_i \|^2 \right]$$

# OBJECTIVE FUNCTION

$$\mathcal{J} = \left[ \sum_{i=1}^{k} \sum_{\mathbf{x} \in C_i} \| P_C^T V^T \mathbf{x} - P_C^T V^T \boldsymbol{\mu}_i \|^2 \right] + \sum_{\mathbf{x} \in \mathcal{D}} \| P_N^T V^T \mathbf{x} - P_N^T \boldsymbol{\mu}_{\mathcal{D}} \|^2$$

# OBJECTIVE FUNCTION

$$\mathcal{J} = \left[ \sum_{i=1}^{k} \sum_{\mathbf{x} \in C_i} \| P_C^T V^T \mathbf{x} - P_C^T V^T \boldsymbol{\mu}_i \|^2 \right]$$
$$+ \sum_{\mathbf{x} \in \mathcal{D}} \| P_N^T V^T \mathbf{x} - P_N^T \boldsymbol{\mu}_{\mathcal{D}} \|^2$$

$P_C \equiv$ cluster space projection matrix

# OBJECTIVE FUNCTION

$$\mathcal{J} = \left[ \sum_{i=1}^{k} \sum_{\mathbf{x} \in C_i} \| P_C^T V^T \mathbf{x} - P_C^T V^T \boldsymbol{\mu}_i \|^2 \right]$$
$$+ \sum_{\mathbf{x} \in \mathcal{D}} \| P_N^T V^T \mathbf{x} - P_N^T \boldsymbol{\mu}_{\mathcal{D}} \|^2$$

$P_C \equiv$ cluster space projection matrix

$P_N \equiv$ noise space projection matrix

# OBJECTIVE FUNCTION

$$\mathcal{J} = \left[ \sum_{i=1}^{k} \sum_{\mathbf{x} \in C_i} \| P_C^T V^T \mathbf{x} - P_C^T V^T \boldsymbol{\mu}_i \|^2 \right]$$
$$+ \sum_{\mathbf{x} \in \mathcal{D}} \| P_N^T V^T \mathbf{x} - P_N^T \boldsymbol{\mu}_{\mathcal{D}} \|^2$$

$P_C \equiv$ cluster space projection matrix

$P_N \equiv$ noise space projection matrix

$V \equiv$ transformation matrix

# OBJECTIVE FUNCTION

$$\mathcal{J} = \mathrm{tr}\left( P_C P_C^T V^T \underbrace{\left( \left[ \sum_{i=1}^{k} S_i \right] - S_{\mathcal{D}} \right)}_{\Sigma} V \right)$$

$$+ \underbrace{\mathrm{tr}(V^T S_{\mathcal{D}} V)}_{\text{const. w.r.t } V}$$

# OBJECTIVE FUNCTION

$$
\mathcal{J} = \mathrm{tr}\left( P_C P_C^T V^T \underbrace{\left( \left[ \sum_{i=1}^{k} S_i \right] - S_{\mathcal{D}} \right)}_{\Sigma} V \right)
$$

$$
+ \underbrace{\mathrm{tr}(V^T S_{\mathcal{D}} V)}_{\mathrm{const.\ w.r.t}\ V}
$$

$$
S_i \equiv \text{cluster scatter matrix}
$$

# OBJECTIVE FUNCTION

$$\mathcal{J} = \mathrm{tr}\left( P_C P_C^T V^T \underbrace{\left( \left[ \sum_{i=1}^{k} S_i \right] - S_{\mathcal{D}} \right)}_{\Sigma} V \right)$$

$$+ \underbrace{\mathrm{tr}(V^T S_{\mathcal{D}} V)}_{\text{const. w.r.t } V}$$

$$S_i \equiv \text{cluster scatter matrix}$$

$$S_{\mathcal{D}} \equiv \text{dataset scatter matrix}$$

# MINIMIZATION

$$\mathcal{J} = \mathrm{tr}\left( P_C P_C^T V^T \underbrace{\left( \left[ \textstyle\sum_{i=1}^{k} S_i \right] - S_{\mathcal{D}} \right)}_{\Sigma} V \right) \ldots$$

- put eigenvectors of $\Sigma$ into $V$ in ascending order

- keep the negative eigenvalues via $P_C P_C^T$

# COMPUTATIONAL COMPLEXITY

$$\mathcal{O}\big(\ I\ (\ mk|\mathcal{D}| + d^2|\mathcal{D}| + d^3\ )\big)$$

# COMPUTATIONAL COMPLEXITY

$$\mathcal{O}\big( \, I \, ( \, mk|\mathcal{D}| + d^2|\mathcal{D}| + d^3 \, ))$$

k-means

# Computational Complexity

$$\mathcal{O}\big(\, I\,(\, mk|\mathcal{D}| + d^2|\mathcal{D}| + d^3\,)\big)$$

k-means

scatter matrix

# COMPUTATIONAL COMPLEXITY

$$\mathcal{O}\big(\ I\ (\ mk|\mathcal{D}| + d^2|\mathcal{D}| + d^3\ )\big)$$

k-means

scatter matrix

eigenvalue decomposition

# Randomized Subspace k-means

# Transformation Matrix Approximation

- $P_C P_C^T$ only keeps the first $m$ eigenvalues

- compute rank-$m$ approximation, $\widetilde{V}$, using a randomized eigenvalue decomposition

# TRANSFORMATION MATRIX APPROXIMATION

- $P_C P_C^T$ only keeps the first $m$ eigenvalues

- compute rank-$m$ approximation, $\widetilde{V}$, using a randomized eigenvalue decomposition

$$\mathcal{J} = \mathrm{tr}\left( \widetilde{V}^T \left( \left[ \sum_{i=1}^{k} S_i \right] - S_{\mathcal{D}} \right) \widetilde{V} \right) + \underbrace{\mathrm{tr}(\widetilde{V}^T S_{\mathcal{D}} \widetilde{V})}_{\text{const. w.r.t } \widetilde{V}}.$$

$\underbrace{\phantom{\widetilde{V}^T \left( \left[ \sum_{i=1}^{k} S_i \right] - S_{\mathcal{D}} \right) \widetilde{V}}}_{\Sigma}$

# COMPUTATIONAL COMPLEXITY

## before:

$$\mathcal{O}\big(I(mk|\mathcal{D}| + d^2|\mathcal{D}| + d^3\,)\big)$$

# COMPUTATIONAL COMPLEXITY

**before:**

$$\mathcal{O}\big(I(mk|\mathcal{D}| + d^2|\mathcal{D}| + d^3\,)\big)$$

**after:**

$$\mathcal{O}\big(I(mk|\mathcal{D}| + d^2|\mathcal{D}| + d^2\log m\,)\big)$$

# COMPUTATIONAL COMPLEXITY

**before:**

$$\mathcal{O}\big(I(mk|\mathcal{D}| + d^2|\mathcal{D}| + d^3 )\big)$$

**after:**

$$\mathcal{O}\big(I(mk|\mathcal{D}| + d^2|\mathcal{D}| + d^2 \log m )\big)$$

# EXPERIMENTS

## Synthetic Data

- runtime vs dimensions
- runtime vs instances

## Real Datasets

- clustering quality
- runtime

## Algorithms

- Subspace k-means
- Randomized Subspace k-means
- PCA k-means
- LDA k-means

# DATASETS

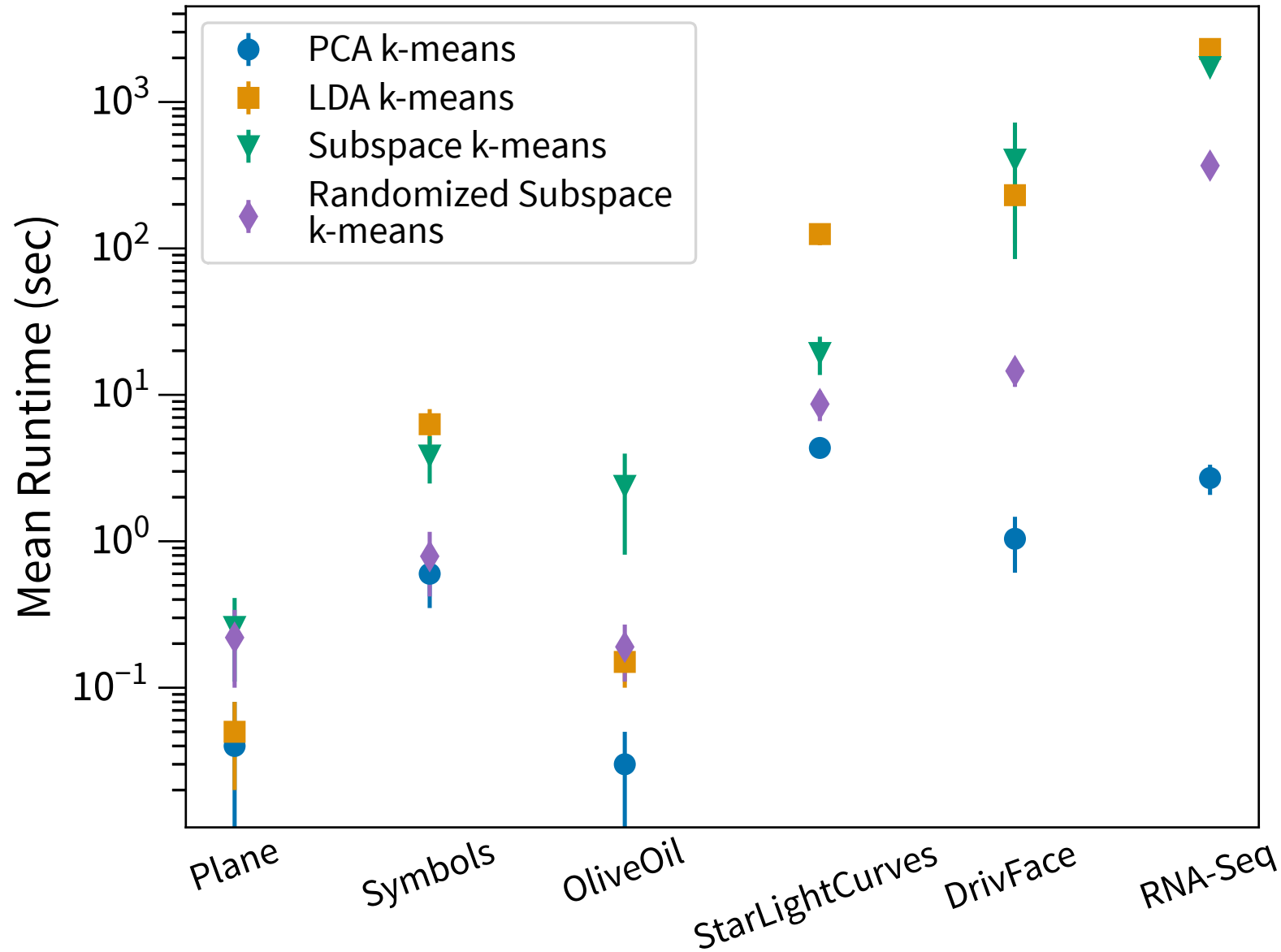| | Features | Instances | Classes |
|---|---|---|---|
| **Plane** | 114 | 210 | 7 |
| **Symbols** | 398 | 1020 | 6 |
| **OliveOil** | 570 | 60 | 4 |
| **StarLightCurves** | 1024 | 9236 | 3 |
| **DrivFace** | 6400 | 606 | 3 |
| **RNA-Seq** | 20531 | 801 | 5 |

# RESULTS

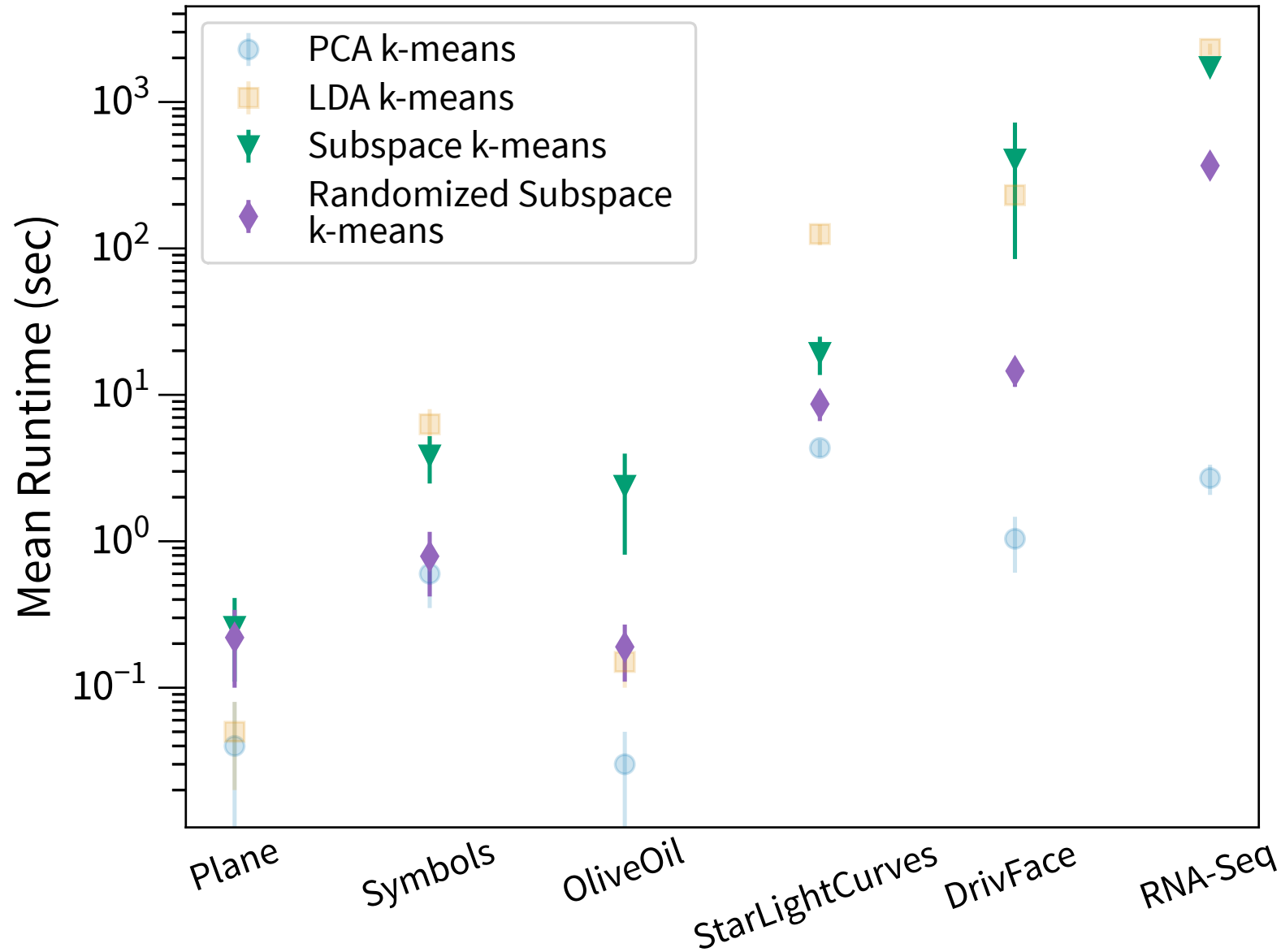**Runtime vs Dimension**

Runtime vs Dimension

# Runtime vs Size

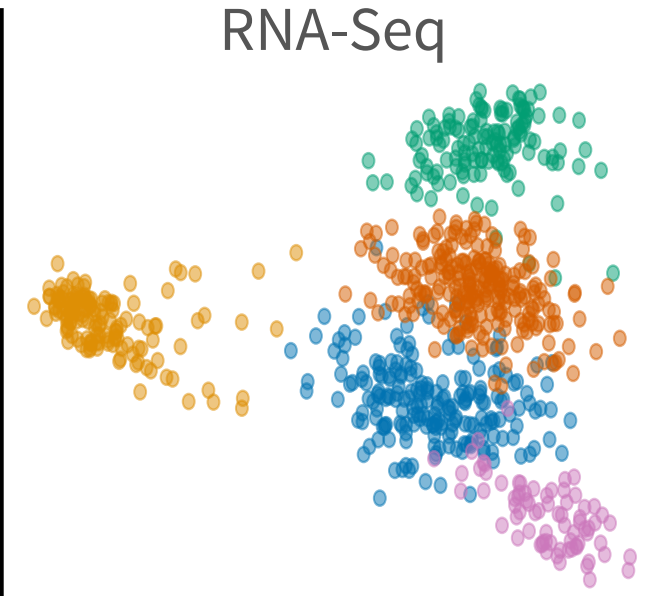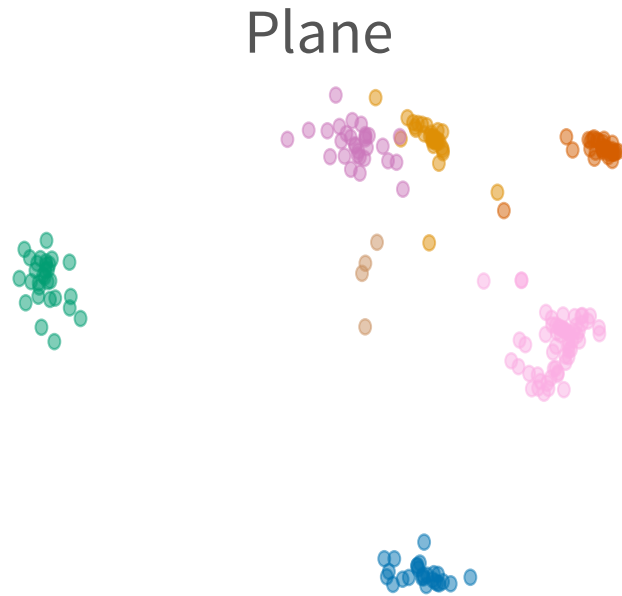# Runtime vs Size

DATASET RUNTIMES

**DATASET RUNTIMES**

# Clustering Quality (NMI)

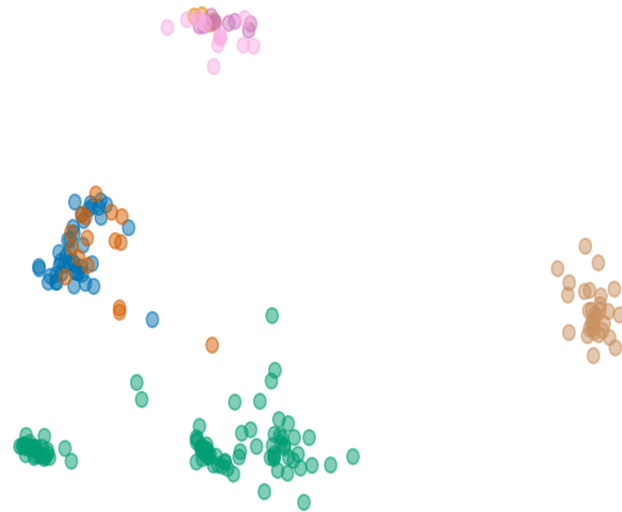| | Randomized Subspace k-means | Subspace k-means | PCA k-means | LDA k-means |
|---|---|---|---|---|
| Plane | **0.835** | 0.825 | 0.804 | 0.728 |
| Symbols | **0.788** | 0.742 | 0.762 | 0.745 |
| OliveOil | 0.609 | 0.657 | 0.673 | **0.690** |
| StarLight Curves | **0.546** | 0.542 | 0.507 | 0.542 |
| DrivFace | 0.191 | 0.205 | 0.203 | **0.209** |
| RNA-Seq | 0.659 | 0.679 | **0.680** | 0.668 |

# SUBSPACE PROJECTIONS

|  | Plane | RNA-Seq |
|---|---|---|
| Two most important features found by *randomized subspace k-means* | | |
| Two most important features found by PCA | | |

# CONCLUSIONS

# Highlights

- significant performance increase
- no reduction in clustering quality

# Future Work

- improve scatter matrix complexity
- k-means extensions
- test on more/larger datasets

QUESTIONS?

# NMI

$$NMI(C,T) = \frac{I(C,T)}{\sqrt{H(C)H(T)}}$$

$C \equiv$ cluster assignments

$T \equiv$ ground truth

$I(C,T) \equiv$ mutual information

$H(C) \equiv$ entropy of cluster assignments

$H(T) \equiv$ entropy of ground truth

# Randomized EVD

Approximate range: $Y = A\Omega$

Obtain orthonormal basis: $Y = QR$

Factorize: $A \approx QQ^*A$

EVD on $B = Q^*A$