#### 3<sup>rd</sup> International Conference On Computational Intelligence And Networks

#### **PRINCIPAL SUBSPACE UPDATION FOR INTEGRATIVE CLUSTERING OF**

#### **MULTIMODAL OMICS DATA**

28 OCTOBER 2017

APARAJITA KHAN and PRADIPTA MAJI MACHINE INTELLIGENCE UNIT INDIAN STATISTICAL INSTITUTE KOLKATA, INDIA.

## Outline

## > Multimodal Data Integration

- ➢ Existing Approaches
- ➢ Principal Subspace
- ► Relevance and Concordance
- ➢ Proposed Algorithm
- ➢ Results and Survival Analysis
- ➤ Conclusion

## Multimodal Data Integration



## Existing Approaches

#### **2 Stage Approaches:**



## Existing Approaches



- Relevance of individual modalities
- Concordance between modalities







## **Relevance and Concordance**

![](_page_7_Picture_1.jpeg)

## Proposed Algorithm

![](_page_8_Figure_1.jpeg)

Computational complexity is  $O(Mn^2d_{max})$ .

- 2 real-life cancer data sets from TCGA
- Modalities:
  - Gene Expression, DNA methylation, miRNA, Protein, CNV
- Giloblastoma Multiforme (GBM):
  - 168 samples, 4534 features, 4 clusters
- Kidney cancer (KIDNEY):
  - 737 samples, 5979 features, 3 clusters
- Compared with:
  - BCC [2], COCA [3]
  - o iCluster [4], LRAcluster [5], PCA-NI
- External Indices: F-measure, NMI
- Survival Analysis: p-value in log-rank test [7]

RESULTS

#### Relevance of Each Modality and Selected Modalities

Kidney

Different	Kidney		GBM	
Modality	Relevance $\mathcal{R}_l$	Selected	Relevance $\mathcal{R}_l$	Selected
CNV	0.2552738		0.2371261	
DNA	0.3840580	Gene	-	Gene
Gene	0.5172488	miRNA	0.2958170	CNV
miRNA	0.4014777	DNA	0.2196443	
Protein	0.2430872		-	

KIDNEY ordering: Gene>miRNA>DNA>CNV>Protein

![](_page_10_Figure_5.jpeg)

## **RESULTS II**

![](_page_11_Figure_1.jpeg)

## **RESULTS III**

![](_page_12_Figure_1.jpeg)

## Survival Analysis

#### Kaplan-Meier survival plots with median survival time

**KIDNEY** 

![](_page_13_Figure_3.jpeg)

GBM

## Conclusion

- Constructs low-rank joint subspace from individual subspaces
- > Evaluates relevance and mutual information before naïve integration
- Filters out noisy and inconsistent modalities
- Computationally efficient

#### Future Work

- $\hfill\square$  Dimensionality increases linearly with number of modalities
- Update subspace instead appending
- □ Effective rank estimation techniques

## References

- H. Zha, X. He, C. Ding, H. Simon, and M. Gu, "Spectral relaxation for k-means clustering," in Neural Information Processing Systems, vol. 14, (Vancouver, Canada), pp. 1057 - 1064, 2001.
- 2. E. F. Lock, et al., "Bayesian consensus clustering," Bioinformatics, no. 29(20), pp. 2610–2616, 2013.
- 3. K. A. Hoadley, C. Yau, et al., "Multiplatform analysis of 12 cancer types reveals molecular classification within and across tissues of origin," Cell, vol. 158, pp. 929–944, 2014.
- R. Shen, A. B. Olshen, and M. Ladanyi, "Integrative clustering of multiple genomic data types using joint latent variable model with application to breast and lung cancer subtype analysis," Bioinformatics, no. 25(22), pp. 2906–2912, 2009.
- 5. D. Wu, D. Wang, et al., , "Fast dimension reduction and integrative clustering of multi-omics data using low-rank approximation: Application to cancer molecular classification," BMC Genomics, 2015.
- 6. D. W. Hosmer, S. Lemeshow, and S. May, Applied Survival Analysis: Regression Modeling of Time to Event Data. New York, NY, USA: Wiley-Interscience, 2nd ed., 2008.
- P. J. Rousseeuw, "Silhouettes: A graphical aid to the interpretation and validation of cluster analysis," Journal of Computational and Applied Mathematics, vol. 20, pp. 53 – 65, 1987.

# THANK YOU!

## **Questions?**