



# Multi-omics data integration by deep learning for phenotype prediction

**Blaise Hanczar** 



Laboratoire IBISC Informatique, Bioinformatique et Systèmes Complexes



## Omics data integration

#### Supervised learning setting

- Gene expression
- DNA Methylation
- Copy Number Variation
- Single Nucléotide Variation
- Meta-genomics
- Clinical
- Medical images

٠ Predictive Model

- Diagnosis
- Prognostics
- Treatment response

How to integrate different sources of data ?

Numeric, binary, categorial, count, sequences, images

## Omics data integration

#### Supervised learning setting

- Gene expression
- DNA Methylation
- Copy Number Variation
- Single Nucléotide Variation
- Meta-genomics
- Clinical
- Medical images



How to integrate different sources of data ? Numeric, binary, categorial, count, sequences, images

### With a neural network

### Integration levels



- Preprocessing
- Feature extraction
- Redundancy
- Note fully exploit the complementary nature of the modalities
- Decisions fusion
- Related to ensemble methods
- Need uncorrelated errors

- Shared representation of the data
- Flexible fusion
- Majority of the deep mutli-modal methods

Ramachandram, D., & Taylor, G. W. (2017). IEEE Signal Processing Magazine, 34(6), 96-108.



Sun, D., Wang, M., & Li, A. (2018). A multimodal deep neural network for human breast cancer prognosis prediction by integrating multi-dimensional data. *IEEE/ACM transactions on computational biology and bioinformatics*, *16*(3), 8<sup>1</sup>/<sub>4</sub>1-850.

### Early integration



Data integration by multi-modal autoencoder

Prediction of liver cancer survival from TCGA dataset

Latent space integrating the information from the three data sources

- K-means
- Cox-PH model
- SVM

Chaudhary, K., Poirion, O. B., Lu, L., & Garmire, L. X. (2018). Deep learning–based multi-omics integration robustly predicts survival in liver cancer. *Clinical Cancer Research*, *24*(6), 1248-1259.



Vale-Silva, L. A., & Rohr, K. (2020). MultiSurv: Long-term cancer survival prediction using multimodal deep learning. *medRxiv*.



Cheerla, A., & Gevaert, O. (2019). Deep learning with multimodal representation for pancancer prognosis prediction. *Bioinformatics*, *35*(14), i446-i454.

8

## Multi omics and knowledge integration



Ma, T., & Zhang, A. (2019). Integrate multi-omics data with biological interaction networks using Multi-view Factorization AutoEncoder (MAE). *BMC genomics*, *20*(11), 1-11.

9

## Some results

Limitations of multi omics models:

- More data sources does not necessarily mean better accuracy
- More data sources implies more costly models

Included data modalities							
Clinical	mRNA	DNAm	miRNA	CNV	ISW	C <sup>td</sup> (95% CI)	IBS (95% CI)
•	•					0.822	0.138
•	•					(0.805-0.837)	(0.126-0.150)
•		•				0.808	0.134
•		•				(0.791-0.826)	(0.125-0.148)
			•			0.792	0.147
•			•			(0.775-0.810)	(0.136-0.161)
•				•		0.795	0.140
•				•		(0.778-0.812)	(0.131-0.152)
•					•	0.801	0.148
•					•	(0.783-0.817)	(0.140-0.158)
•	•	•				0.810	0.146
•	•	•				(0.793 - 0.829)	(0.135-0.158)
•	•	•	•			0.798	0.153
•	•	•	•			(0.781-0.815)	(0.139-0.168)
•		•	•	•		0.802	0.149
•	•	•	•	•		(0.748 - 0.820)	(0.136-0.162)
	•		•	•	•	0.787	0.152
•	•	•	•	•	•	(0.769 - 0.806)	(0.140-0.166)

Clinical: tabular clinical data; mRNA: gene expression; DNAm: DNA methylation; miRNA: microRNA expression; CNV: gene copy1number variation; WSI: whole-slide images.

## Which data sources to use ?



 $Cost = c_1 + c_2 + c_3$ 

- Tradeoff between accuracy and cost
- Some patients are easy to predict :
  - Do no need all data sources
  - Prediction could be less expensive with the same accuracy

## Budget learning for multi omics data

- Objectives in classification:
  - Maximize the accuracy
  - Minimize the prediction cost
- Cost is mainly the cost of acquisition of the variables
  - Money
  - Time
  - Secondary effect
  - Any non-infinite ressource
- Adaptative integration of omics data

## Adaptative integration of omics data

- Sequence of neural networks with reject option of increasing cost
- At each level, we compute :
  - Prediction from current data sources
  - Decision to acquire the next data source



 $Cost = c_1 + c_2 + c_3$  13

## End-to-End learning



- $R_1$  and  $R_2$  control the integration of  $x_1$  and  $x_2$
- The order of the data sources is important