### Multi-modal data analysis

- with an application to brain functional genomics

Daifeng Wang, Ph.D. Assistant Professor Department of Biostatistics and Medical Informatics Department of Computer Sciences (affiliate) Investigator, Waisman Center University of Wisconsin – Madison daifeng.wang@wisc.edu https://daifengwanglab.org/ BMI/CS 776 Spring 2022

## Machine learning to open "black box" from genotype to phenotype

- Goal
  - Advance biological knowledge on genomics in brain diseases
- Approach
  - Machine learning, Bioinformatics, Computational Biology

#### Decoding Genomic Information to Better Understand Molecular Mechanisms and Improve Disease Diagnosis



### Outline

- Background
  - Functional genomics & multi-omics
  - Gene regulation
- Multi-view learning for multi-modal data
  - Multi-view Empirical Risk Minimization (MV-ERM)
  - Manifold alignment
- App 1: single-cell multi-modal data integration
  - Predicting neuronal electrophysiology from gene expression
- App 2: *deep manifold-regularized classification* 
  - Predicting cortical layers of neurons from multi-modalities
- App 3: biologically interpretable neural network modeling
  - Disease prediction and functional prioritization

### Your genome is your genetic code book

Book	Genome
Chapters	Chromosomes
Sentences	Genes
Words	Elements
Letters	Bases

#### Human

- 46 chromosomes
- $\sim 20,000 25,000$  genes
- ~ Millions elements
- 4 unique bases (A, T, C, G), ~3 billion in total



### How to read your genetic code book?



Chapter One Republicans and Democrats

	Peek	Conomo
	BOOK	Genome
2	Chapters	Chromosomes
	Sentences	Genes
	Words	Elements
	Letters	Bases

"On most days, I enter the Capitol through the basement. A small subway train carries me from the Hart Building, where ..."

• Key words

#### • Non-key words

Overhead, the ceiling forms a creamy white oval, with an American eagle etched in its center. Above the visitors' gallery, the busts of the nation's first twenty vice presidents sit in solemn repose.

And in gentle steps, one hundred mahogany desks rise from the well of the Senate in four horseshoe-shaped rows. Some of these desks date back to 1819, and atop each desk is a tidy receptacle for inkwells and quills. Open the drawer of any desk, and you will find within the names of the senators who once used it—Taft and Long, Stennis and Kennedy—scratched or penned in the senator's own hand. Sometimes, standing there in



### Grammar for book is clear but not for genome



#### Low sequencing cost enables reading our whole genome



#### Cost per Genome

#### Differences across our genomes: DNA variants

### Single Nucleotide Polymorphisms (SNPs) normally happen ~1% on individual human genome.



Illustration by Darryl Leja, NHGRI

### Genome-Wide Association Study (GWAS) identifies disease associated genetic variants

#### 36,989 schizophrenia cases and 113,075 controls in Psychiatric Genomics Consortium



Schizophrenia Working Group of the Psychiatric Genomics Consortium, Nature (2014)

However, genotype-phenotype association can't tell mechanisms

### How to link non-coding disease SNPs to genes?



🔺 TF 🔶 RBP 🖕 miRNA 🔺 Germ-line variant 🐇 Somatic variant



### Complex mechanisms from genotype to phenotype



Diseaseassociated genomic variants



## Hierarchical understanding from genotype to phenotype



### Multi-omics data for various functional elements





Based on an image by Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

### Gene expression and regulation

Identical DNA but different gene expression





**Gene regulation**: mechanisms controlling gene expression levels

### Gene regulatory networks linking functional elements from multi-omics



Jin et al, Genome Medicine, 2021

### Gene regulatory networks linking GWAS SNPs to disease genes

1. Predicting cell-type gene regulatory networks (GRNs) via multi-omics



#### 3. Microglia GRN, cell-type disease gene functions



### 2. Identifying cell-type disease genes via GWAS and cell-type GRNs





Ting Jin (Ph.D. student, Biomedical Data Science)



## Machine learning for multi-omics analysis, prediction and interpretation



Xu. Genome Bio. 2019

16

### Outline

- Background
  - Functional genomics & multi-omics
  - Gene regulation
- Multi-view learning for multi-modal data
  - Multi-view Empirical Risk Minimization (MV-ERM)
  - Manifold alignment
- App 1: single-cell multi-modal data integration
  - Predicting neuronal electrophysiology from gene expression
- App 2: *deep manifold-regularized classification* 
  - Predicting cortical layers of neurons from multi-modalities
- App 3: biologically interpretable neural network modeling
  - Disease prediction and functional prioritization

#### **Empirical Risk Minimization (ERM)**



## Empirical Risk Minimization (ERM) for single-omics (single-view learning)



- e.g., Leukemia patient classification
  - *y<sub>i</sub>*: Acute lymphoblastic leukemia
     (ALL) vs. Acute myeloid leukemia
     (AML)
  - $x_i$ : gene expression

– f: SVM



### Example: single-view unsupervised learning

## Non-negative matrix factorization (NMF)





Single cell signatures

• ~14,000 cells (Lake et al.,

Science, 2016&2018)

### Multi-view learning for multi-omics integration



- For example, gene regulation involves
  - 1. Genomics; e.g., SNPs
  - 2. Transcriptomics; e.g., genes
  - 3. Proteomics; e.g., transcription factors (TFs)

Cross-omics interactions  $\Omega_{co}(f^{(1)}, f^{(3)})$ : SNPs break TF binding sites  $\Omega_{co}(f^{(2)}, f^{(3)})$ : TFs control gene expression  $\Omega_{co}(f^{(1)}, f^{(2)})$ : SNPs associate with gene expression (e.g, eQTLs)

### Multiview Empirical Risk Minimization (MV-ERM)



Relationship information across views

#### Factorization-based MV-ERM



### Alignment-based MV-ERM



### Alignment-based multi-view learning methods

#### Multiple kernel learning

non-linear

- learn a combination of <u>predefined</u> kernels
- not data dependent
- expensive
  - Only support <u>complementary principle</u>

 $\beta_i K^{(i)}$ 

#### Subspace learning

- obtains a common latent subspace
- **data dependent**  $\rightarrow$  more general
- non-expensive
- CCA (Canonical Correlation Analysis)
  - Only support <u>consensus principle</u>
  - linear  $\max_{f,g} corr(f(X), g(Y))$



Can we support **both principles**?

## ManiNetCluster: manifold alignment to reveal the functional links between gene networks



(by employing **manifold learning**) & support <u>both principles</u>

Nguyen, Blaby, Wang, BMC Genomics, 2019 (Best Poster Award, ACM BCB 2018)

## ManiNetCluster: manifold alignment to reveal the functional links between gene networks

Application: genomic functional linkages between light and dark periods of green alga



functional linkage (Module 34) cell motility UDP glucosyl

### Manifold Learning



A *d* dimensional manifold *M* is embedded in a *m* dimensional space, and there is an explicit mapping  $f: \mathbb{R}^m \to \mathbb{R}^d$  where  $d \leq m$ . Given samples  $x_i \in \mathbb{R}^m$  with noise



### Manifold Alignment

 $\begin{array}{ll} two \quad X = [x_1, \dots, x_m], x_i \in \mathbb{R}^p \\ datasets \quad Y = [y_1, \dots, y_n], y_j \in \mathbb{R}^q \end{array} \quad x_i \leftrightarrow y_i \ for \ i \in [1, l] \end{array}$ 



### Manifold Alignment framed by MV-ERM

 $\begin{aligned} & \text{Learn common} \quad X = [x_1, \dots, x_m], x_i \in \mathbb{R}^p \\ & \text{manifolds between two} \quad Y = [y_1, \dots, y_n], y_j \in \mathbb{R}^q \quad x_i \leftrightarrow y_i \text{ for } i \in [1, l] \\ & \text{To find mapping function } f(.), g(.) \text{ to minimize the cost function} \\ & \text{Inter-view correspondence} \quad & \text{Intra-view correspondence} \\ & \sum_{i,j} \left\| f(x_i) - g(y_j) \right\|^2 W^{i,j} + \sum_{i,j} \left\| f(x_i) - f(x_j) \right\|^2 W_X^{i,j} + \sum_{i,j} \left\| g(y_i) - g(y_j) \right\|^2 W_Y^{i,j} \\ & \text{global consistency} \quad & \text{local smoothness} \\ & \text{complementary} \end{aligned}$ 

- Extract and optimally align local geometry to minimize overall differences
- Generalization of CCA  $\sum_{i,j} \|f(x_i) g(y_j)\|^2$
- Interpreted as a manifold regularization

$$\sum_{i,j} \left\| f(x_i) - g(y_j) \right\|^2 W^{i,j} + tr(\mathbb{f}^T L_X \mathbb{f}) + tr(\mathbb{g}^T L_Y \mathbb{g})$$

• Eigen-decomposition to solve nonlinear manifold alignment

### Outline

- Background
  - Functional genomics & multi-omics
  - Gene regulation
- Multi-view learning for multi-modal data
  - Multi-view Empirical Risk Minimization (MV-ERM)
  - Manifold alignment
- App 1: single-cell multi-modal data integration
  - Predicting neuronal electrophysiology from gene expression
- App 2: deep manifold-regularized classification
  - Predicting cortical layers of neurons from multi-modalities
- App 3: biologically interpretable neural network modeling
  - Disease prediction and functional prioritization

### Single-cell multi-modal data by Patch-seq (beyond omics)



# Manifold alignment of single neurons by electrophysiology and gene expression





Jiawei Huang (M.S., Statistics, 2021) Patch-seq data of ~3k neuronal cells in the mouse visual cortex from BRAIN Initiative

### Nonlinear Manifold Alignment (NMA) uncovers trajectory with multi-modal changes



# NMA aligned cells form cross-modal clusters, suggesting related genes and electrophysiological features



### Predicting electrophysiological features by differentially expressed genes in cross-modal clusters



### Outline

- Background
  - Functional genomics & multi-omics
  - Gene regulation
- Multi-view learning for multi-modal data
  - Multi-view Empirical Risk Minimization (MV-ERM)
  - Manifold alignment
- App 1: single-cell multi-modal data integration
  - Predicting neuronal electrophysiology from gene expression
- App 2: *deep manifold-regularized classification* 
  - Predicting cortical layers of neurons from multi-modalities
- App 3: biologically interpretable neural network modeling
  - Disease prediction and functional prioritization

# *deepManReg*: a deep manifold-regularized learning model for phenotype prediction from multi-modal data



### Phase 1: deep manifold alignment of multi-modal features

- Goal: reveal nonlinear relationships across multi-modal features
- Trade-off:
  - Features  $\rightarrow$  Parametric
  - Nonlinear  $\rightarrow$  Nonparametric
- deepManReg solution:
  - Parameterize  $f(\cdot, \mathcal{W}) \& g(\cdot, \mathcal{Z})$  by deep neural networks
  - Nonlinear manifold alignment (NMA)
- Eigen-decomposition to solve NMA:
  - Computationally intensive
  - Use gradient descent on a non-Euclidean space!



### Optimization for deepManReg alignment

 $\mathbb{F} = [f(X; \mathcal{W})^T, g(Y; \mathcal{Z})^T]^T$  can be solved by

$$\min_{f,g} tr(\mathbb{F}^T L \mathbb{F}) \text{ s.t. } \mathbb{F}^T D \mathbb{F} = \mathbb{I},$$

where *L* and *D* are joint Laplacian and diagonal matrices of prior correspondence across modalities  $\rightarrow$  an optimization problem on Stiefel manifold  $O^{d \times r}$ 



• Forward pass – project the output  $\begin{bmatrix} f(X) \\ g(Y) \end{bmatrix}$  onto Stiefel manifold

$$F = \pi \circ \begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} = UIV^T$$

, where  $\begin{bmatrix} f(X) \\ g(Y) \end{bmatrix} = U\Sigma V^T$  is the SVD of output

- Backward pass
  - Compute Euclidean gradient

$$\nabla_F \ell = \frac{\partial tr(F^T L F)}{\partial F} = LF + L^T F$$

• Project Euclidean gradient onto the tangent space of Stiefel manifold

$$\widetilde{\nabla}_F \ell = \pi(\nabla_F \ell) = Fskew(F^T \nabla_F \ell) + (I - FF^T) \nabla_F \ell$$

 $\circ \quad \text{Backpropagate the Riemannian gradient } \widetilde{\nabla}_F \ell \text{ to} \\ \text{update } \mathcal{W}, \mathcal{Z}$ 

Cunningham et al. The Journal of Machine Learning Research 16.1 (2015): 2859-2900.

# Phase 2: classification regularized by cross-modal feature network (aligned features)



## deepManReg alignment & classification of handwritten digits



#### Multi-modal feature alignment

- 2000 images of handwritten digits 0-9 (mfeat data) Breukelen et al. Kybernetika, 34(4):381–386, 1998.
- Two types of features:
  216 profile correlations, 76 Fourier coefficients



## deepManReg alignment & classification of neuronal cells

Single-cell multi-modal data by Patch-seq

- 3654 GABAergic cortical neurons in mouse visual cortex
- <u>Modality 1</u>: Electrophysiological features
  - By hyperpolarizing and depolarizing current injection stimuli and responses of short (3 ms) current pulses, long (1 s) current steps, and slow (25 pA/s) current ramps.
- <u>Modality 2</u>: Genes
  - Expression levels
- <u>Phenotype</u>: 5 cortical layers, L1, L2/3, L4, L5, and L6
- Running time (alignment only)
  - CCA (725.96 seconds)
  - Manifold Alignment (663.43 seconds) MATCHER (150.94 seconds)
  - deepManReg (57.90 seconds by GPUs GTX 1060Ti and 90.10 seconds by CPU i5-8250U)



# Cross-modal feature network linking aligned genes and electrophysiological features

Electrophysiological feature

Gene



e-Features	Index
rheobase_sweep_num	1
thumbnail_sweep_num	2
vrest	3
sag	4
tau	5
f_i_curve_slope	6
adaptation	7
latency	8
upstroke_downstroke_ratio_long_square	9
peak_v_long_square	10
peak_t_long_square	11
trough_v_long_square	12
trough_t_long_square	13
fast_trough_v_long_square	14
fast_trough_t_long_square	15
threshold_i_long_square	16
threshold_t_long_square	17
upstroke_downstroke_ratio_ramp	18
peak_v_ramp	19
peak_t_ramp	20
trough_v_ramp	21
trough_t_ramp	22
fast_trough_v_ramp	23
fast_trough_t_ramp	24
threshold_t_ramp	25
upstroke_downstroke_ratio_short_square	26
peak_v_short_square	27
peak_t_short_square	28
trough_t_short_square	29
fast_trough_t_short_square	30
threshold_v_short_square	31
threshold_t_short_square	32

## Feature-network-regularized classification of cortical layers of neuronal cells



### Outline

- Background
  - Functional genomics & multi-omics
  - Gene regulation
- Multi-view learning for multi-modal data
  - Multi-view Empirical Risk Minimization (MV-ERM)
  - Manifold alignment
- App 1: single-cell multi-modal data integration
  - Predicting neuronal electrophysiology from gene expression
- App 2: *deep manifold-regularized classification* 
  - Predicting cortical layers of neurons from multi-modalities
- App 3: biologically interpretable neural network modeling
  - Disease prediction and functional prioritization

## Machine learning opens "black box" for disease prediction and functional interpretation



# *Varmole*: Interpretable deep neural network model prioritizes disease variants and genes via DropConnect



- Input form 2 omics, X, Y (SNPs & genes)
- First layer embed  $A_1$  and  $A_2$  eQTL and gene regulatory network (GRN)

 $\rightarrow$  From variants (& gene regulations) to gene expression

• Other fully connected hidden layers: *h*;

 $\rightarrow$  From gene expression to phenotypes

- Softmax classification layer:  $o = \delta \left( h \circ \sigma (f(X) + g(Y)) \right);$
- The Cross-Entropy:  $L(o, \hat{o}) = -\frac{1}{n \sum_{i=1}^{n} y_i log(\hat{y}_i)}$
- Varmole: min  $L(o, \hat{o}) + ||W||_1$

#### **Drop-connect**

• Drop-out and drop-connect are 2 simple but effective regularization techniques



The drop-connect mask is eQTL or GRN ( $A_1$  or  $A_2$ )

### **Interpretation:** prioritization via Integrated Gradients

- Given a model *F*, an input x, and the output F(x) of the model for input in question, an attribution methods returns the 'relevance' of each input feature *i* to the output
- How importance a gene (SNP, eQTL link, TF-gene link) to the disease outcome



 $\frac{\partial F(x)}{\partial x_i}$  is the gradient of F along the  $i^{th}$  dimension at x

Importance score

of link  $x_i \rightarrow y$ 

### Application for schizophrenia

- Dataset:
  - RNA-seq gene expression & genotype data (dosage) for 487 schizophrenia (scz) vs. 891 non-scz human brain samples (front cortex)
  - Embedding GTEx eQTLs & PsychENCODE GRN for human brain front cortex
  - $\rightarrow$  127304 SNPs, 2598 genes



## Prioritized gene functions & regulatory links for schizophrenia

- A list of enriched functions (FDR<0.05) from prioritized genes:
  - neuron development
  - axon guidance
  - cell adhesion
  - calcium signaling
  - response to external stimulus
  - NMDA receptor
  - insulin secretion
- Prioritized SNP-gene pairs
  - SNP-gene pairs on the interacting enhancers and promoters (Hi-C) have significantly higher importance scores (*p*<5e-5)</li>
  - Potential regulatory roles of prioritized SNPs to genes via enhancers

### Outline

- Background
  - Functional genomics & multi-omics
  - Gene regulation
- Multi-view learning for multi-modal data
  - Multi-view Empirical Risk Minimization (MV-ERM)
  - Manifold alignment
- App 1: single-cell multi-modal data integration
  - Predicting neuronal electrophysiology from gene expression
- App 2: *deep manifold-regularized classification* 
  - Predicting cortical layers of neurons from multi-modalities
- App 3: biologically interpretable neural network modeling
  - Disease prediction and functional prioritization

### Thank you!

#### Current lab members

- Dr. Shuang Liu (Data Scientist III)
- Dr. Pramod Chandrashekar (Postdoc)
- Dr. Chirag Gupta (Postdoc)
- Dr. Chenfeng He (Postdoc)
- Dr. Kalpana H. Arachchilage (Postdoc)

- Saniya Khullar (CIBM fellow & Ph.D. candidate, Biomedical Data Science)
- Ting Jin (Ph.D. student, Biomedical Data Science)
- Sayali Alatkar (Ph.D. student, Computer Sciences)
- Noah Cohen Kalafut (Ph.D. student, Computer Sciences)
- Jie Sheng (Biostatistician, Waisman Center)
- Jonathan Bryan (Undergraduate, Neurobiology)

#### Recent graduates

- Nam Nguyen (Ph.D. 2021, Computer Science, Lane Fellow in Carnegie Mellon University)
- Jiawei Huang (M.S. 2021, Statistics, Ph.D. student in University of Cincinnati)
- Yudi Mu (M.S. 2021, Statistics, Ph.D. student in University of California Riverside)
- Mufang Yin (M.S. 2020, Statistics, Ph.D. student in Rutgers University)
- Peter Rehani (B.S. 2020, Neurobiology, Morgridge Institute for Research)

#### Funding acknowledgment

- NSF Career 2144475
- NIH R01AG067025
- NIH R21CA237955
- NIH R21NS127432
- NIH R03NS123969
- NIH U01MH116492
- NIH P50HD105353
- NIH R01HD106197
- DOE QPSI

<u>daifeng.wang@wisc.edu</u> <u>https://daifengwanglab.org/</u> *Various positions available*!

