# Machine Learning Network-Constrained Regression of Epigenetic Data

Sivo Vladimirov Daskalov

Corpus Christi College

28 June 2017

UNIVERSITY OF
CAMBRIDGE

# Outline

# Epigenetic background



Gene "switched on"
- Active (open) chromatin
- Unmethylated cytosines (white circles)
- Acetylated histones

Transcription Factors / Co-activators

Gene "switched off"
- Silent (condensed) chromatin
- Methylated cytosines (red circles)
- Deacetylated histones

Transcription possible

SWI/SNF

HAT

RNA Pol II

HDAC

HMT

Transcription impeded

*Figure is adapted from Luong, P. Basic Principles of Genetics

# Project goals

Question:
How is the expression of each gene affected by the methylation of related genes?

Approach:

Linear regression $\begin{cases} \text{Predictors: methylation levels for all genes} \\ \text{Target variable: expression level for gene of interest} \end{cases}$

# Penalized regression methods

$$\text{Lasso } \lambda \sum_{i=1}^{p} |\beta_i|$$

$$\text{Elastic Net } \lambda_1 \sum_{i=1}^{p} |\beta_i| + \lambda_2 \sqrt{\sum_{i=1}^{p} \beta_i^2}$$

$$\text{Grace } \lambda_1 \sum_{i=1}^{p} |\beta_i| + \lambda_2 \sum_{u \sim v} \left( \frac{\beta_u}{\sqrt{d_u}} - \frac{\beta_v}{\sqrt{d_v}} \right)^2 w(u, v)$$

$$\text{aGrace } \lambda_1 \sum_{i=1}^{p} |\beta_i| + \lambda_2 \sum_{u \sim v} \left( \frac{sign(\tilde{\beta}_u)\beta_u}{\sqrt{d_u}} - \frac{sign(\tilde{\beta}_v)\beta_v}{\sqrt{d_v}} \right)^2 w(u, v)$$

$$\text{GBLasso } \lambda \sum_{u \sim v} \left[ \left( \frac{|\beta_u|}{\sqrt{d_u}} \right)^{\gamma} + \left( \frac{|\beta_v|}{\sqrt{d_v}} \right)^{\gamma} \right]^{1/\gamma}$$

$$\text{Linf } \lambda \sum_{u \sim v} \max \left( \frac{|\beta_u|}{\sqrt{d_u}}, \frac{|\beta_v|}{\sqrt{d_v}} \right)$$

$$\text{aLinf } \lambda \sum_{u \sim v} \left| \frac{sign(\tilde{\beta}_u)\beta_u}{\sqrt{d_u}} - \frac{sign(\tilde{\beta}_v)\beta_v}{\sqrt{d_v}} \right|$$

$$\text{TTLP } \lambda_1 \sum_{i=1}^{p} J_\tau |\beta_i| + \lambda_2 \sum_{u \sim v} \left| J_\tau \left( \frac{|\beta_u|}{w_u} \right) - J_\tau \left( \frac{|\beta_v|}{w_v} \right) \right|$$

$$\text{LTLP } \lambda_1 \sum_{i=1}^{p} |\beta_i| + \lambda_2 \sum_{u \sim v} \left| J_\tau \left( \frac{|\beta_u|}{w_u} \right) - J_\tau \left( \frac{|\beta_v|}{w_v} \right) \right|$$
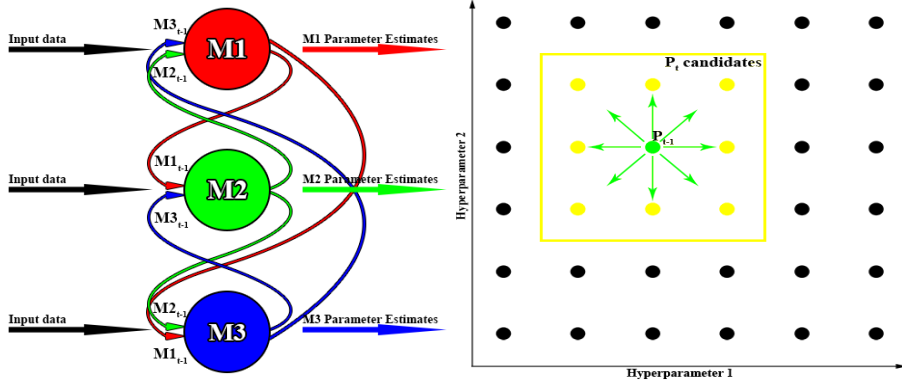
# Composite voting regression

|          | $X_1$         | $X_2$         | ...  | $X_p$         |
|----------|---------------|---------------|------|---------------|
| Method 1 | $M_1(\beta_1)$ | $M_1(\beta_2)$ | ...  | $M_1(\beta_p)$ |
| Method 2 | $M_2(\beta_1)$ | $M_2(\beta_2)$ | ...  | $M_2(\beta_p)$ |
| ...      | ...           | ...           | ...  | ...           |
| Method k | $M_k(\beta_1)$ | $M_k(\beta_2)$ | ...  | $M_k(\beta_p)$ |

$$X_j = \begin{cases} important, & \text{if } \frac{\sum_{i=1}^{k}[M_i(\beta_j) \neq 0]}{k} \geq \text{fraction of votes threshold} \\ unrelated, & \text{otherwise} \end{cases}$$

Final model obtained from OLSE on the set of important predictors

# Orchestrated hyperparameter tuning

# Synthetic dataset generation and setup

Synthetic dataset generation

- ▶ Designed to be similar to real epigenetic datasets
- ▶ 20 datasets with 550 predictors and differently generated responses
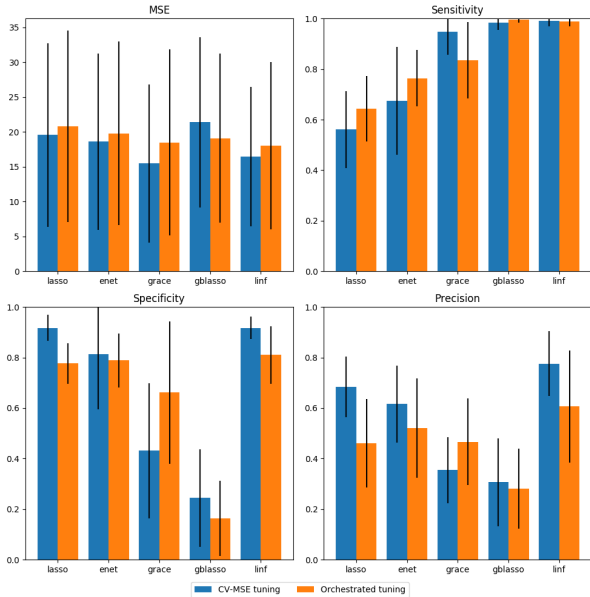- ▶ Training and test sets of size 300 and 100 respectively

Hyperparameter tuning setup

Search space: Predefined parameter grids for all regression methods

CV-MSE tuning: Traditional 5-fold cross-validated mean squared error

Orchestrated tuning: Starting points obtained from the CV-MSE tuning

# Comparison of model metrics

# Regression method similarity evaluation

Cosine similarity between estimated coefficient vectors



CV-MSE tuning

Orchestrated tuning

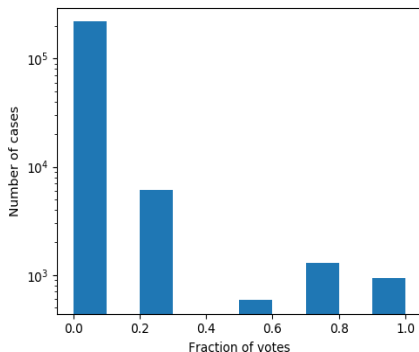# Breast cancer dataset and setup

Dataset properties:
- ▶ Obtained from The Cancer Genome Atlas (TCGA)
- ▶ Methylation and expression data for 215 breast cancer patients
- ▶ Selected subset of genes associated with breast cancer
- ▶ Samples divided in 3/4 training and 1/4 test sets

Methylation data from the promoter and gene body regions considered separately
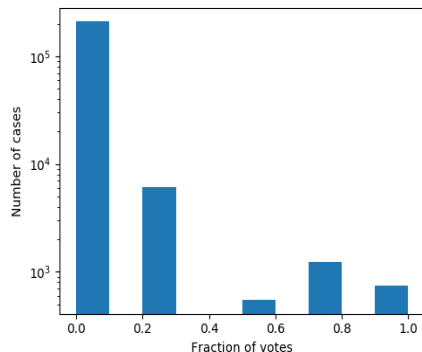
Regression methods used: Lasso, Elastic Net, Grace, Linf and the proposed Composite Voting Regression

# Vote fraction distribution

Threshold of 0.75 chosen (3 out of 4 methods must agree)
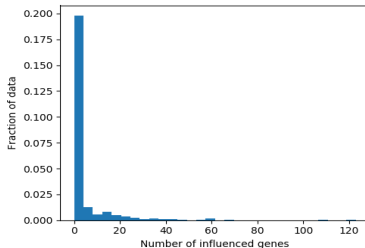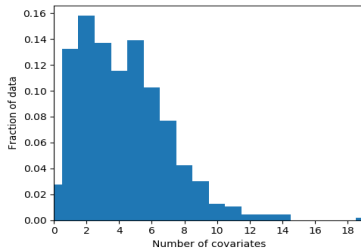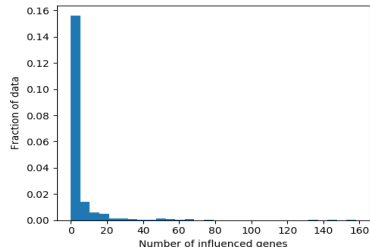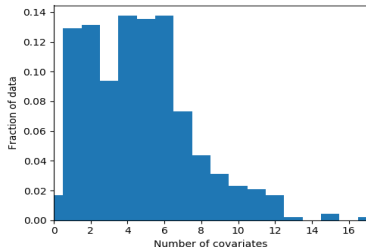


Gene body region

Gene promoter region

# Distribution of dependencies

## Composite voting regression on the gene promoter region



## Composite voting regression on the gene body region

# Summary

- Implementation of 9 regression methods found in literature
- Composite voting regression
- Orchestrated hyperparameter tuning
- Comparison and evaluation on synthetic datasets
- Exploration of a real breast cancer dataset

Contact details:
sivodaskalov@gmail.com