

# HiC Data Visualization and Statistical Analysis

Amal Agarwal Advisors: Dr. Yu Zhang and Dr. Lingzhou Xue

> Department of Statistics Pennstate University

> > May 25, 2018





2 Methodology



2 Methodology



#### HiC data structure:

- Two Cell types, "Gm12878" and "K562".
- 22 chromosomes and X chromosome in each cell type.
- Intensity matrix for each chromosome with 10K b.p. granularity.

## Challenge:

Intensity matrix is big!





#### HiC data structure:

- Two Cell types, "Gm12878" and "K562".
- 22 chromosomes and X chromosome in each cell type.
- Intensity matrix for each chromosome with 10K b.p. granularity.

## Challenge:

Intensity matrix is big!





2 Methodology



# Regression Model



#### **Notation:**

- N: size of HiC matrix for one chromosome.
- $\blacksquare$   $n_b$  observations (banded values) for bandwidth parameter b.
- Consider only upper triangular elements (i.e.  $j \ge i$ ).

## Penalized Lasso Regression Model:

Intensity<sub>ij</sub> =  $\alpha + \beta \times |i - j| + \sum_{k=1}^{N} \gamma_k \mathcal{I}(i \le k \le j) + \epsilon$ subject to  $\beta \le 0, \gamma_k \le 0, \ \forall k \ \text{and} \ \sum_{k=1}^{N} \gamma_k \ge t$ , where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ 





#### **Notation:**

- N: size of HiC matrix for one chromosome.
- $\blacksquare$   $n_b$  observations (banded values) for bandwidth parameter b.
- Consider only upper triangular elements (i.e.  $j \ge i$ ).

### Penalized Lasso Regression Model:

Intensity<sub>ij</sub> =  $\alpha + \beta \times |i - j| + \sum_{k=1}^{N} \gamma_k \mathcal{I}(i \le k \le j) + \epsilon$ subject to  $\beta \le 0, \gamma_k \le 0, \ \forall k \ \text{and} \ \sum_{k=1}^{N} \gamma_k \ge t$ , where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ 



# Regression Model



# Equivalent formulation as Penalized Least Squares (PLS):

$$\min \left( \sum_{i=1}^{N} \sum_{j=i}^{\min(N,i+n_b)} \left[ \left( \text{Intensity}_{ij} - \alpha - \beta - \sum_{k=1}^{N} \gamma_k \mathcal{I}(i \leq k \leq j) \right)^2 - \lambda \left( \sum_{k=1}^{N} \gamma_k \right) \right] \right)$$

#### Challenges due to large N:

- For  $N \sim 25 K$  and  $b \sim 200$ , we have  $n_b \sim 10 M$ .
- The design matrix becomes  $\sim 10M \times 25K \implies 25B$  elements. Memory problem!
- Solving the PLS function is computationally expensive





## Equivalent formulation as Penalized Least Squares (PLS):

$$\min \left( \sum_{i=1}^{N} \sum_{j=i}^{\min(N,i+n_b)} \left[ \left( \text{Intensity}_{ij} - \alpha - \beta - \sum_{k=1}^{N} \gamma_k \mathcal{I}(i \leq k \leq j) \right)^2 - \lambda \left( \sum_{k=1}^{N} \gamma_k \right) \right] \right)$$

#### Challenges due to large N:

- For  $N \sim 25 K$  and  $b \sim 200$ , we have  $n_b \sim 10 M$ .
- The design matrix becomes  $\sim 10M \times 25K \implies 25B$  elements. Memory problem!
- Solving the PLS function is computationally expensive.





## Coordinate Descent (CD) Algorithms:

- Used by "glmnet"
- Convergence is fast
- Closed form updates
- Active set strategy





## Objective (generalized notation):

$$f_{\lambda}(\alpha, \beta, \gamma) = \sum_{i=1}^{n} \left[ (y_i - \alpha - \beta z_i - \gamma^T x_i)^2 - \lambda \sum_{j=1}^{p} \gamma_j \right]$$

$$(\hat{\alpha}, \hat{\beta}, \hat{\gamma}) = \arg\min (f_{\lambda}(\alpha, \beta, \gamma))$$
  $\beta \leq 0,$   $\gamma \leq 0$ 

**Note:** Tuning parameter  $\lambda$  is assumed to be known here.



# Naive CD Algorithm



#### Algorithm 1: Naive Coordinate Descent

**Input:** Tuning Parameter  $\lambda$ , Convergence threshold  $\epsilon$ 

Output:  $\hat{\alpha}, \hat{\beta}, \hat{\gamma}$ 

- 1 Initialization: Start with arbitrary  $(lpha^{(1)},eta^{(1)},\gamma^{(1)})$
- 2 Compute  $f_{\lambda}^{(1)} = f_{\lambda}(\alpha^{(1)}, \beta^{(1)}, \gamma^{(1)})$
- 3 Set  $f_1^{(0)} \leftarrow \infty$  and  $t \leftarrow 1$
- 4 while  $|f_{\lambda}^{(t)} f_{\lambda}^{(t-1)}| > \epsilon$  do

5 
$$\alpha^{(t+1)} \leftarrow \arg\min_{\alpha} (f_{\lambda}(\alpha, \beta^{(t)}, \gamma^{(t)}))$$

6 
$$\beta^{(t+1)} \leftarrow \arg\min_{\beta} (f_{\lambda}(\alpha^{(t+1)}, \beta, \gamma^{(t)}))$$

7 for 
$$j \leftarrow 1$$
 to  $p$  do

Compute 
$$f_{\lambda}^{(t+1)} = f_{\lambda}(\alpha^{(t+1)}, \beta^{(t+1)}, \gamma^{(t+1)})$$
 and set  $t \leftarrow t+1$ 





## Check out the HiC report for:

- Active Set coordinate descent algorithm.
- Details about the closed form updates.





2 Methodology

**Future Directions** 



### In process...

- Coding the Coordinate Descent algorithms.
- Exploring other stopping criteria.

## Questions:

■ How to design simulations that mimick HiC data?

