

Distribution Learning ppt

Wei Ying, Yiming Li

November 29, 2023

Overview

1 Introduction

2 Method

3 Simulation

Problem Statement

Suppose in the reference site, we observe n i.i.d. samples (X_i, T_i, \mathbf{G}_i) , where the genotypes of the i -th subject on a segment of DNA of interest $\mathbf{G}_i \in \mathbb{R}^p$, the gene expression level of the i -th subject in the RNA sequence data $X_i \in \mathbb{R}^1$, and the phenotypes of the i -th subject in a GWAS data $T_i \in \mathbb{R}^1$.

Problem Statement

Assuming conditional independence of \mathbf{T} and $\mathbf{X}_{\mathbf{G}}$ given \mathbf{G} , we have

$$\begin{aligned}\mathbb{E}[t|x] &= \frac{\iint t f(t|x, \mathbf{G}) f(x|\mathbf{G}) f(\mathbf{G}) d\mathbf{G} dt}{\int_{\mathbf{G}} f(x|\mathbf{G}) f(\mathbf{G}) d\mathbf{G}} \\ &= \frac{\iint t f(x|\mathbf{G}) f(t, \mathbf{G}) d\mathbf{G} dt}{\int_{\mathbf{G}} f(x|\mathbf{G}) f(\mathbf{G}) d\mathbf{G}} \\ &\approx \frac{\sum_i T_i f(\mathbf{X}_{\mathbf{G}} = x | G_i)}{\sum_i f(\mathbf{X}_{\mathbf{G}} = x | G_i)}\end{aligned}$$

The key step in our proposed approach is to estimate the conditional density of $\mathbf{X}_{\mathbf{G}}|\mathbf{G}$.

Dimension Reduction

Since \mathbf{G} is high-dimension, we consider the principal component to catch the main feature of \mathbf{G} .

$$\operatorname{argmin}_{\mathbf{A}, \mathbf{B}, f} \text{Cost} = \|\mathbf{G} - \mathbf{A}\mathbf{B}^\top\| + L(\mathbf{X}, f(\mathbf{A})) + \Omega(\mathbf{A}, \mathbf{B})$$

where \mathbf{B} defines the feature of SNPs and \mathbf{A} is individual mutation loading on that factor or the projection of genotype data on these factors. The first part is the decomposition of the aimed matrix \mathbf{G} . The second part is the measurement of the similarity of our proposed distribution $f(\mathbf{A})$ to the real \mathbf{X} .

A natural way is to consider the Gaussian mixture for the conditional distribution of $\mathbf{X} \mid \mathbf{G}$. If the genotypes can be approximated by D features, i.e., $\mathbf{A}_{n \times D}$, it makes sense to assume the gene expression is also a mixture of D Gaussian distributions, whose means and variances depend on $\mathbf{A}_d, d = 1, 2, \dots, D$. So we could assume

$$\mathbf{X} \mid \mathbf{G} \sim \sum_{d=1}^D \pi_d \Psi(\mu_d(\mathbf{A}_d), \sigma_d(\mathbf{A}_d))$$

In theory, any continuous function can be approximated by the combination of basis functions with a sufficient number of components. We hope $\mu_d(A_d)$ can be well approximated by a combination of K basis functions:

$$\mathbf{A}_d = (A_{1d}, A_{2d}, \dots, A_{nd}), \boldsymbol{\mu}_d = (\mu_{1d}, \mu_{2d}, \dots, \mu_{nd}), \boldsymbol{\sigma}_d = (\sigma_{1d}, \sigma_{2d}, \dots, \sigma_{nd})$$

$$\mu_{id} = \sum_{q=1}^K \alpha_q B_q(A_{id})$$

We also hope σ_{id} is the MLE estimation, i.e., $\sigma_{id} = (\sum_{i=1}^n (X_i - \mu_{id}))^{\frac{1}{2}}$

One Dimension Case

When \mathbf{X} is a n by 1 matrix. The Optimization question is:

$$C = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^p (g_{ij} - a_i b_j)^2 + \sum_{i=1}^n \log \sigma_i + \frac{(x_i - \mu_i)^2}{2\sigma_i^2}$$

$$\text{subject to } \|B_1\|_2 = \sqrt{\sum_j^p b_j^2} = 1$$

$$\text{where } \sigma_i = \left(\sum_{i=1}^n (X_i - \mu_i) \right)^{\frac{1}{2}}, \quad \mu_i = \sum_{q=1}^K \alpha_q B_q(a_i)$$

One Dimension Case

$$a_i = \frac{\sum_{j=1}^p g_{ij} b_j - \frac{(\mu_i - x_i) \mu'_i}{\sigma_i^2}}{\sum_{j=1}^p b_j^2}, \quad b_j = \frac{\sum_{i=1}^n a_i g_{ij}}{\sum_{i=1}^n a_i^2}, \quad \alpha = E\left[\frac{\mathbf{X}}{\sigma} \middle| \frac{B_{q=1:K}}{\sigma}\right]$$

B should be sparse, we consider the standard Lasso penalization on b_j and get the following solutions:

$$b_j^{lasso} = \text{sgn}(b_j^*) \left(|b_j^*| - \frac{\lambda}{\sum_{i=1}^n a_i^2} \right)^+, \quad \text{where } b_j^* = \frac{\sum_{i=1}^n a_i g_{ij}}{\sum_{i=1}^n a_i^2}$$

Additive Model

Since the genotype data is high-dimensional, the computation time would be hugely prolonged if we always refresh the previous updated value and rebuild the model. One possible solution might be residual regression. In our case, We fix the previously updated coefficients, \mathbf{A}_{pre} , \mathbf{B}_{pre} , α_{pre} where pre represents previous updated value. Then we only update the residual of \mathbf{G} : $\mathbf{G}_{\text{res}} = \mathbf{G}_{\text{pre}} - \mathbf{A}_{\text{pre}} \mathbf{B}_{\text{pre}}^{\top}$ where res represents residual. To include as much information on genotype data \mathbf{G} as possible, we require \mathbf{B}_{new} to be orthogonal with all the previous \mathbf{B}_{pre} . We realize this by multiplying \mathbf{B}_{new} by a matrix $\mathbf{I}_p - \mathbf{B}_{\text{pre}} \left(\mathbf{B}_{\text{pre}}^{\top} \mathbf{B}_{\text{pre}} \right)^{-1} \mathbf{B}_{\text{pre}}^{\top}$ to project \mathbf{B}_{new} on the orthogonal space of \mathbf{B}_{pre} .

Algorithm for One Dimension Case

S1: Start with $D = 1$

- s1 update a_i, b_j, α_q for $i = 1:n, j = 1:p, q = 1:K$
- s2 calculate Γ matrix for Gaussian mixture model
- s3 compare where the 2-norms of difference for $\mathbf{A}, \mathbf{B}, \alpha$ are smaller than the prespecified value. If yes, calculate the cost and move forward to **Step 2** $D = D+1$. If not, go back to s1

S2: $\mathbf{G} = \mathbf{G} - \mathbf{AB}^\top$

- ss1 update a_i, b_j, α_q for $i = 1:n, j = 1:p, q = 1:K$. Multiply \mathbf{B}_{new} by $\mathbf{I}_p - \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$.
- ss2 calculate Γ matrix for Gaussian mixture model
- ss3 compare where the 2-norms of difference for $\mathbf{A}, \mathbf{B}, \alpha$ are smaller than the prespecified value. If not, go back to ss1. If yes, calculate the cost and move forward to ss7.
- ss4 If new cost decrease more than 10%, repeat **Step 2**, $D = D+1$. If not, Stop

One Dimension With Multiple Response

When \mathbf{X} is n by M , $M > 1$

$$C = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^p (g_{ij} - a_i b_j)^2 + \sum_{m=1}^M \left(\sum_{i=1}^n \log \sigma_i^m + \frac{(x_i^m - \mu_i^m - \boldsymbol{\eta}_m \mathbf{x}_i^{-m})^2}{2(\sigma_i^m)^2} \right)$$

$$\text{subject to } \|B_1\|_2 = \sqrt{\sum_j b_j^2} = 1$$

$$\text{where } \sigma_i^m = \left(\sum_{i=1}^n (x_i^m - \mu_i^m) \right)^{\frac{1}{2}}, \quad \mu_i^m = \sum_{q=1}^K \alpha_q^m B_q(a_i)$$

$$\boldsymbol{\eta}_m = (\eta_m^1, \dots, \eta_m^{m-1}, \eta_m^{m+1}, \eta_m^M)$$

$$a_i^m = \frac{\sum_{j=1}^p g_{ij} b_j - \sum_{m=1}^M \frac{(\mu_i^m + \eta_m \mathbf{x}_i^{-m} - \mathbf{x}_i^m) \mu_i^{m'}}{\sigma_i^{m2}}}{\sum_{j=1}^p b_j^2}$$

$$b_j^m = \frac{\sum_{i=1}^n a_i g_{ij}}{\sum_{i=1}^n a_i^2}$$

$$\alpha^m = E\left[\frac{\mathbf{X}^m - \eta_m \mathbf{X}^{-m}}{\sigma^m} \middle| \frac{B_{q=1:K}}{\sigma^m}\right]$$

Solution

$$\Delta = \begin{pmatrix} 1 & \eta_{1,1}^{(d^*)} & \cdots & \eta_{1,m-2}^{(d^*)} & \eta_{1,m-1}^{(d^*)} \\ \eta_{2,1}^{(d^*)} & 1 & \cdots & \eta_{2,m-2}^{(d^*)} & \eta_{2,m-1}^{(d^*)} \\ & & \cdots & & \\ \eta_{m,1}^{(d^*)} & \eta_{m,2}^{(d^*)} & \cdots & \eta_{m,m-1}^{(d^*)} & 1 \end{pmatrix}$$

Δ is symmetric about the diagonal. Thus we have extra constraints when solving η .

Algorithm for One Dimension Case With Multiple Response

S1: Start with $D = 1$

- s1 update $a_i^m, b_j^m, \alpha_q^m, \boldsymbol{\eta}^m$ for $i = 1:n, j = 1:p, q = 1:K, m = 1:M$
- s2 calculate Γ matrix for Gaussian mixture model
- s3 compare where the 2-norms of difference for $\mathbf{A}, \mathbf{B}, \boldsymbol{\alpha}$ are smaller than the prespecified value. If yes, calculate the cost and move forward to **Step 2** $D = D+1$. If not, go back to s1

S2: $\mathbf{G} = \mathbf{G} - \mathbf{AB}^\top$

- ss1 update $a_i^m, b_j^m, \alpha_q^m, \boldsymbol{\eta}^m$ for $i = 1:n, j = 1:p, q = 1:K, m = 1:M$.
Multiply \mathbf{B}_{new} by $\mathbf{I}_p - \mathbf{B}(\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$.
- ss2 calculate Γ matrix for Gaussian mixture model
- ss3 compare where the 2-norms of difference for $\mathbf{A}, \mathbf{B}, \boldsymbol{\alpha}$ are smaller than the prespecified value. If not, go back to ss1. If yes, calculate the cost and move forward to ss7.
- ss4 If new cost decrease more than 10%, repeat **Step 2**, $D = D+1$. If not, Stop

Simulation Study

Aim1: one dimension case, our proposed method is better than neural network and glm. WD qq plot.

Aim2: if we have multiple responses \mathbf{X} , We believe \mathbf{X} are inner correlated. Thus, Conditional model(our) should be better than Marginal model(duplicate one response model).

Simulation Result

We use the whole blood sample data GTE_xXS_{NHG5}. The number of individuals is 670, the number of snp is 3417. We set the causal.rate = 0.01. For the homogeneous model:

$$Y = \beta^T \mathbf{X}_{\text{casual}} + \epsilon$$

For the heterogeneous model:

$$Y = \beta^T \mathbf{X}_{\text{casual}} + (1 + c\beta^T \mathbf{X}_{\text{casual}})\epsilon$$

We then normalize Y and remove the highly duplicate column X.

Single Result

We use 10 fold cross-validation to test the performance of three methods. WD is the abbreviation for Wasserstein distance.

$$W_1(\mu_1, \mu_2) = \int_{\mathbb{R}} |F_1(y) - F_2(y)| dx$$

We use numerical approximation:

$$W_1(\mu_1, \mu_2) = \frac{1}{n} \sum_{i=1}^n |Q_1(y, \frac{i}{n}) - Q_2(y, \frac{i}{n})|$$

Where $Q(y, \frac{i}{n})$ is the $\frac{i}{n}$ -th quantile of y . WD is the result for our proposed method; WD2 is the result for the elastic net with optimal tuning parameter; WD3 is the result for neural network method with the same number of mixtures as our proposed method.

Single Result

For this homo case, we set the coefficient $\beta = 0.1$

```
```{r}
abs()
mean(WD)
mean(WD2)
mean(WD3)

var(WD)
var(WD2)
var(WD3)
```
```

[1] 0.1723963
[1] 0.4081771
[1] 0.3144389
[1] 0.003727351
[1] 0.0004586147
[1] 0.03686145

Figure: homosingle

Single Result

For this hete case, we set the coefficient $\beta = 0.1$

```
```{r}
abs()
mean(WD)
mean(WD2)
mean(WD3)

var(WD)
var(WD2)
var(WD3)
```

[1] 0.1723963
[1] 0.4081771
[1] 0.3144389
[1] 0.003727351
[1] 0.0004586147
[1] 0.03686145
```

Figure: hetesingle

Multiple Result

We also want to show when \mathbf{Y} 's are correlated with each other, the multiple-algorithm could work better than multiple single-algorithms.

$$\mathbf{Y} = \beta^\top \mathbf{X} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \Sigma)$$

The result of multiple-algorithm are simulated using the Markov chain

$$Y_i^t \sim \mathcal{N}(\mu + \eta_i \mathbf{Y}_{-i}^{(t-1)}, \sigma_i) \quad i = 1, 2, \dots, D \quad t = 1, 2, \dots, T$$

Where i represents the position of dimension \mathbf{Y} , and t represents the time-varying variable. And \mathbf{Y}^0 is some random initial values. The prediction is the mean of the latter sequence of $\{\mathbf{Y}^t\} \mid t \geq 0.75T$

cm is the abbreviation for the MSE of the conditional model (multiple-version), mm is the abbreviation for the MSE of the marginal model (single-version). cmwd is the abbreviation for the Wasserstein distance of the conditional model (multiple-version), mmwd is the abbreviation for the Wasserstein distance of the marginal model (single-version).

Multiple Result

For homo case, Σ is the identity matrix. Conditional model vs marginal model (mse)

```
```{r}
cm1
cm2
cm3

(cm1 + cm2 + cm3)/3
```
```

```
[1] 0.831271
[1] 0.6225814
[1] 0.5336311
[1] 0.6624945
```

```
```{r}
mm1
mm2
mm3

(mm1 + mm2 + mm3)/3
```
```

```
[1] 0.8257151
[1] 0.8897963
[1] 0.8620893
[1] 0.8592002
```

Multiple Result

Conditional model vs marginal model (Wasserstein distance)

```
```{r}  
cmwd1
cmwd2
cmwd3

(cmwd1 + cmwd2 + cmwd3)/3
```
```

```
[1] 0.4922632  
[1] 0.2026474  
[1] 0.1487395  
[1] 0.2812167
```

```
```{r}  
mmwd1
mmwd2
mmwd3

(mmwd1 + mmwd2 + mmwd3)/3
```
```

```
[1] 0.4563343  
[1] 0.2794005  
[1] 0.3870225  
[1] 0.3742525
```

Figure: homo

Multiple Result

For hete model,

$$\mathbf{Y} = \boldsymbol{\beta}^\top \mathbf{X} + (1 + c\boldsymbol{\beta}^\top \mathbf{X})\epsilon$$

We suppose $\boldsymbol{\beta}^\top = (0.1, 0.3, 0.5)$ and $\Sigma = \begin{pmatrix} 3 & 2 & 1 \\ 2 & 3 & 2 \\ 1 & 2 & 3 \end{pmatrix}$ and $D = 3$.

1,2,3 represents the position of dimension.

Multiple Result

conditional model vs marginal model (mse)

```
## {r}  
cm1  
cm2  
cm3  
  
(cm1 + cm2 + cm3 )/3  
##
```

```
[1] 0.8223781  
[1] 0.7486099  
[1] 0.7678286  
[1] 0.7796055
```

```
## {r}  
mm1  
mm2  
mm3  
  
(mm1 + mm2 + mm3 )/3  
##
```

```
[1] 0.7801134  
[1] 0.7755223  
[1] 0.8705353  
[1] 0.8087237
```

Figure: hete

Multiple Result

conditional model vs marginal model (Wasserstein distance)

```
```{r}
cmwd1
cmwd2
cmwd3

(cmwd1 + cmwd2 + cmwd3)/3
```
```

```
[1] 0.4737388
[1] 0.4178665
[1] 0.3876963
[1] 0.4264339
```

```
```{r}
mmwd1
mmwd2
mmwd3

(mmwd1 + mmwd2 + mmwd3)/3|
```
```

```
[1] 0.5586712
[1] 0.3910544
[1] 0.4188093
[1] 0.4561783
```

Figure: hete