

Machine Learning Applications in Genetics and Genomics

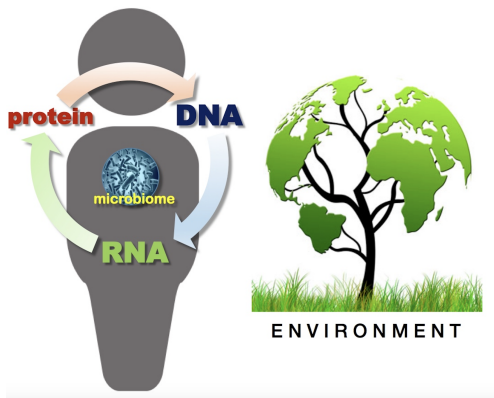
Jing Ma

Public Health Sciences Division
Fred Hutch Cancer Research Center
jingma@fredhutch.org

SLAB LAB
February 5, 2018

Part I:
Improved Variance Component Score Tests of
Gene-Environment Interaction

Gene-environment interaction



Gene-environment interaction

$$y_i = \alpha + \beta_1 \underbrace{E_i}_{\text{environment}} + \beta'_2 \underbrace{G_i}_{\text{genetic markers}} + \gamma' \underbrace{G_i E_i}_{\text{interaction}} + \varepsilon_i,$$

Gene-environment interaction

$$y_i = \alpha + \beta_1 \underbrace{E_i}_{\text{environment}} + \beta'_2 \underbrace{G_i}_{\text{genetic markers}} + \gamma' \underbrace{G_i E_i}_{\text{interaction}} + \varepsilon_i,$$

- y : health outcome (continuous)

Gene-environment interaction

$$y_i = \alpha + \beta_1 \underbrace{E_i}_{\text{environment}} + \beta'_2 \underbrace{G_i}_{\text{genetic markers}} + \gamma' \underbrace{G_i E_i}_{\text{interaction}} + \varepsilon_i,$$

- ▶ y : health outcome (continuous)
- ▶ $\varepsilon_i \sim N(0, \sigma^2)$

Gene-environment interaction

$$y_i = \alpha + \beta_1 \underbrace{E_i}_{\text{environment}} + \beta_2' \underbrace{G_i}_{\text{genetic markers}} + \gamma' \underbrace{G_i E_i}_{\text{interaction}} + \varepsilon_i,$$

- ▶ y : health outcome (continuous)
- ▶ $\varepsilon_i \sim N(0, \sigma^2)$
- ▶ Want to test $H_0 : \gamma = 0$

Connection with linear mixed models

$$y_i = \alpha + \beta_1 E_i + \beta_2' G_i + \underbrace{\gamma'}_{\text{random effects}} G_i E_i + \varepsilon_i,$$

¹ Liu et al. Biometrics. 2007

Connection with linear mixed models

$$y_i = \alpha + \beta_1 E_i + \beta_2' G_i + \underbrace{\gamma'}_{\text{random effects}} G_i E_i + \varepsilon_i,$$

- **Exercise:** find the connection between least squares and linear mixed models¹.

¹ Liu et al. Biometrics. 2007

Connection with linear mixed models

$$y_i = \alpha + \beta_1 E_i + \beta_2' G_i + \underbrace{\gamma'}_{\text{random effects}} G_i E_i + \varepsilon_i,$$

- ▶ **Exercise:** find the connection between least squares and linear mixed models¹.
- ▶ Random effects $\gamma_j \sim N(0, \tau)$.

¹ Liu et al. Biometrics. 2007

Connection with linear mixed models

$$y_i = \alpha + \beta_1 E_i + \beta_2' G_i + \underbrace{\gamma'}_{\text{random effects}} G_i E_i + \varepsilon_i,$$

- ▶ **Exercise:** find the connection between least squares and linear mixed models¹.
- ▶ Random effects $\gamma_j \sim N(0, \tau)$.
- ▶ No interaction $H_0 : \gamma = 0 \Leftrightarrow H_0 : \tau = 0$.

¹ Liu et al. Biometrics. 2007

Variance component test

The VC score statistic is

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

Variance component test

The VC score statistic is

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ \hat{y}_0 is the predicted mean of y .

Variance component test

The VC score statistic is

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ \hat{y}_0 is the predicted mean of y .
- ▶ K depends on G and E .

Variance component test

The VC score statistic is

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ \hat{y}_0 is the predicted mean of y .
- ▶ K depends on G and E .
- ▶ $\lambda_1, \dots, \lambda_n$ are the eigenvalues of $\sigma^2 P_0 K P_0$.

Variance component test

The VC score statistic is

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ \hat{y}_0 is the predicted mean of y .
- ▶ K depends on G and E .
- ▶ $\lambda_1, \dots, \lambda_n$ are the eigenvalues of $\sigma^2 P_0 K P_0$.
- ▶ $P_0 = \mathbf{I}_n - H$ where H is the hat (projection) matrix.

Estimation

Null model:

$$\begin{aligned}y &= \alpha + \beta_1 E + G\beta_2 + \varepsilon \\ &= X\beta + \varepsilon\end{aligned}$$

- ▶ β can be estimated via OLS or ridge regression.

Estimation

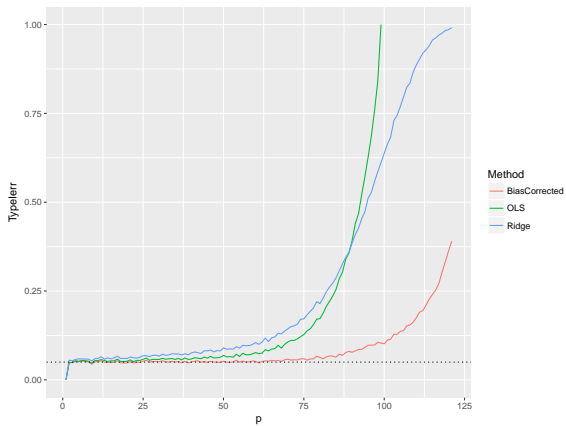
Null model:

$$\begin{aligned}y &= \alpha + \beta_1 E + G\beta_2 + \varepsilon \\ &= X\beta + \varepsilon\end{aligned}$$

- ▶ β can be estimated via OLS or ridge regression.
- ▶ σ^2 can be estimated from RSS.

Type I error - simulation

- $n = 100, \beta_1 = 1, \beta_{2,1} = 1, \sigma^2 = 1.$



Why OLS and Ridge fail

- ▶ OLS does not do well when p is large.

² Lin et al. Biometrics. 2016

Why OLS and Ridge fail

- ▶ OLS does not do well when p is large.
- ▶ The ridge estimator² of β is biased!

$$\hat{\beta} = \arg \min \left\{ \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right\}$$

² Lin et al. Biometrics. 2016

Why OLS and Ridge fail

- ▶ OLS does not do well when p is large.
- ▶ The ridge estimator² of β is biased!

$$\hat{\beta} = \arg \min \left\{ \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right\}$$

- ▶ The bias is not critical when all signals are weak, but causes inflated type I errors when some signals are strong.

² Lin et al. Biometrics. 2016

Ridge estimator is biased

The ridge regression

$$\begin{aligned}\hat{\beta} &= \arg \min \left\{ \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right\} \\ &= (X'X + \lambda I_p)^{-1} X'y\end{aligned}$$

Ridge estimator is biased

The ridge regression

$$\begin{aligned}\hat{\beta} &= \arg \min \left\{ \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 \right\} \\ &= (X'X + \lambda I_p)^{-1} X'y\end{aligned}$$

Let $P_X = X'(XX')^{-1}X$ and $\theta^0 = P_X\beta^0$. The bias in $\hat{\beta}_j$ is

$$\underbrace{\mathbb{E}[\hat{\beta}_j] - \theta_j^0}_{\text{estimation bias}} + \underbrace{\theta_j^0 - \beta_j^0}_{\text{projection bias}}$$

Bias corrected ridge estimator

- ▶ The estimation bias is dominated by the choice of λ .

³ Bühlmann. Bernoulli. 2013

Bias corrected ridge estimator

- ▶ The estimation bias is dominated by the choice of λ .
- ▶ The projection bias

$$B_j = (P_X \beta^0)_j - \beta_j^0 = (P_X)_{jj} \beta_j^0 - \beta_j^0 + \sum_{k \neq j} (P_X)_{jk} \beta_k^0.$$

³ Buhlmann. Bernoulli. 2013

Bias corrected ridge estimator

- ▶ The estimation bias is dominated by the choice of λ .
- ▶ The projection bias

$$B_j = (P_X \beta^0)_j - \beta_j^0 = (P_X)_{jj} \beta_j^0 - \beta_j^0 + \sum_{k \neq j} (P_X)_{jk} \beta_k^0.$$

- ▶ The bias corrected ridge estimator³ is

$$\hat{\beta}_j^{corr} = \hat{\beta}_j - \sum_{k \neq j} (P_X)_{jk} \hat{\beta}_k^{init}.$$

³ Buhlmann. Bernoulli. 2013

Bias corrected ridge estimator

- ▶ The estimation bias is dominated by the choice of λ .
- ▶ The projection bias

$$B_j = (P_X \beta^0)_j - \beta_j^0 = (P_X)_{jj} \beta_j^0 - \beta_j^0 + \sum_{k \neq j} (P_X)_{jk} \beta_k^0.$$

- ▶ The bias corrected ridge estimator³ is

$$\hat{\beta}_j^{corr} = \hat{\beta}_j - \sum_{k \neq j} (P_X)_{jk} \hat{\beta}_k^{init}.$$

- ▶ $\hat{\beta}^{init}$ can be from (scaled) lasso.

³ Bühlmann. Bernoulli. 2013

Property of the bias corrected estimator

$$\hat{\beta}_j^{corr} = Z_j + \delta_j, \quad j = 1, \dots, p$$

- ▶ Z_j 's are normal.

Property of the bias corrected estimator

$$\hat{\beta}_j^{corr} = Z_j + \delta_j, \quad j = 1, \dots, p$$

- ▶ Z_j 's are normal.
- ▶ $\delta_j = (P_X)_{jj}\beta_j^0 - \sum_{k \neq j} (P_X)_{jk}(\hat{\beta}_k^{init} - \beta_k^0) + \mathbb{E}[\hat{\beta}_j] - \theta_j^0$.

Property of the bias corrected estimator

$$\hat{\beta}_j^{corr} = Z_j + \delta_j, \quad j = 1, \dots, p$$

- ▶ Z_j 's are normal.
- ▶ $\delta_j = (P_X)_{jj}\beta_j^0 - \sum_{k \neq j} (P_X)_{jk}(\hat{\beta}_k^{init} - \beta_k^0) + \mathbb{E}[\hat{\beta}_j] - \theta_j^0$.
- ▶ $\hat{\beta}^{corr}$ is easy to compute.

The matrix P_0

- ▶ $P_0 = \mathbf{I}_n - X(X'X + \lambda \mathbf{I}_p)^{-1} X'$ when the estimator is $\hat{\beta}^R$.

The matrix P_0

- ▶ $P_0 = \mathbf{I}_n - X(X'X + \lambda \mathbf{I}_p)^{-1} X'$ when the estimator is $\hat{\beta}^R$.
- ▶ What is P_0 if using $\hat{\beta}^{corr}$?

$$\begin{aligned}\hat{\beta}^{corr} &= \hat{\beta}^R - \{P_X - \text{diag}(P_X)\} \hat{\beta}^{init} \\ &= \hat{\beta}^R - \{P_X - \text{diag}(P_X)\} \text{diag}(W) \hat{\beta}^R \\ &= [\mathbf{I}_p - \{P_X - \text{diag}(P_X)\} \text{diag}(W)] \hat{\beta}^R\end{aligned}$$

The matrix P_0

- ▶ $P_0 = \mathbf{I}_n - X(X'X + \lambda \mathbf{I}_p)^{-1} X'$ when the estimator is $\hat{\beta}^R$.
- ▶ What is P_0 if using $\hat{\beta}^{corr}$?

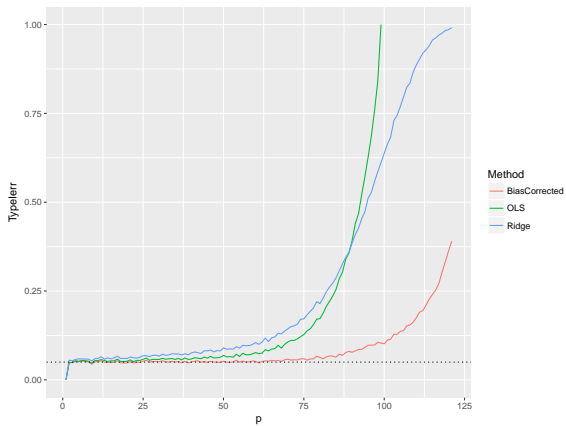
$$\begin{aligned}\hat{\beta}^{corr} &= \hat{\beta}^R - \{P_X - \text{diag}(P_X)\} \hat{\beta}^{init} \\ &= \hat{\beta}^R - \{P_X - \text{diag}(P_X)\} \text{diag}(W) \hat{\beta}^R \\ &= [\mathbf{I}_p - \{P_X - \text{diag}(P_X)\} \text{diag}(W)] \hat{\beta}^R\end{aligned}$$

- ▶ New P_0 with $\hat{\beta}^{corr}$ is

$$\mathbf{I}_n - X[\mathbf{I}_p - \{P_X - \text{diag}(P_X)\} \text{diag}(W)](X'X + \lambda \mathbf{I}_p)^{-1} X'.$$

Type I error - simulation

- $n = 100, \beta_1 = 1, \beta_{2,1} = 1, \sigma^2 = 1.$



Open questions

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ Other bias correction procedures that allow valid type I error control for $p \sim n$, and possibly $p > n$.

Open questions

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ Other bias correction procedures that allow valid type I error control for $p \sim n$, and possibly $p > n$.
 - ▶ Improved fit of the residual $y - \hat{y}_0$.

Open questions

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ Other bias correction procedures that allow valid type I error control for $p \sim n$, and possibly $p > n$.
 - ▶ Improved fit of the residual $y - \hat{y}_0$.
 - ▶ Or even better, estimate directly $X'(y - \hat{y}_0)$.

Open questions

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ Other bias correction procedures that allow valid type I error control for $p \sim n$, and possibly $p > n$.
 - ▶ Improved fit of the residual $y - \hat{y}_0$.
 - ▶ Or even better, estimate directly $X'(y - \hat{y}_0)$.
- ▶ What if y is binary?

Open questions

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ Other bias correction procedures that allow valid type I error control for $p \sim n$, and possibly $p > n$.
 - ▶ Improved fit of the residual $y - \hat{y}_0$.
 - ▶ Or even better, estimate directly $X'(y - \hat{y}_0)$.
- ▶ What if y is binary?
- ▶ Kernel machine based interaction testing.

Open questions

$$Q = (y - \hat{y}_0)' K (y - \hat{y}_0) \sim \sum_{i=1}^n \lambda_i \chi_{1,i}^2$$

- ▶ Other bias correction procedures that allow valid type I error control for $p \sim n$, and possibly $p > n$.
 - ▶ Improved fit of the residual $y - \hat{y}_0$.
 - ▶ Or even better, estimate directly $X'(y - \hat{y}_0)$.
- ▶ What if y is binary?
- ▶ Kernel machine based interaction testing.
- ▶ Beyond VC score test.

Another look at the model

Suppose E_i is binary (e.g. smoker / non-smoker).

$$\begin{aligned}y_i &= \alpha + \beta_1 E_i + \beta_2' G_i + \gamma' G_i E_i + \varepsilon_i, \\&= \alpha + \beta_1 E_i + (\beta_2 + \gamma E_i)' G_i \\&= \begin{cases} \alpha + \beta_2' G_i, & \text{if } E_i = 0 \\ \alpha + \beta_1 + (\beta_2 + \gamma)' G_i, & \text{if } E_i = 1. \end{cases}\end{aligned}$$

Two sample inference of HD linear regressions

$$y = \alpha_d + X\beta^{(d)} + \varepsilon_d, \quad d = 1, 2$$

Two sample inference of HD linear regressions

$$y = \alpha_d + X\beta^{(d)} + \varepsilon_d, \quad d = 1, 2$$

► $H_{0,j} : \beta_j^{(1)} = \beta_j^{(2)}, j = 1, \dots, p.$

Two sample inference of HD linear regressions

$$y = \alpha_d + X\beta^{(d)} + \varepsilon_d, \quad d = 1, 2$$

- ▶ $H_{0,j} : \beta_j^{(1)} = \beta_j^{(2)}, j = 1, \dots, p.$
- ▶ Inverse regression X_j on (y, X_{-j}) (Xia et al. Stat. Sin. 2017)

Two sample inference of HD linear regressions

$$y = \alpha_d + X\beta^{(d)} + \varepsilon_d, \quad d = 1, 2$$

- ▶ $H_{0,j} : \beta_j^{(1)} = \beta_j^{(2)}, j = 1, \dots, p.$
- ▶ Inverse regression X_j on (y, X_{-j}) (Xia et al. Stat. Sin. 2017)
- ▶ Alternative: direct comparison of $\beta_j^{(1)}$ and $\beta_j^{(2)}$?

From linear to logistic regressions

- ▶ First approach:

$$\text{logit}P(y_i = 1) = \alpha_d + X_i' \beta^{(d)}, \quad d = 1, 2$$

From linear to logistic regressions

- ▶ First approach:

$$\text{logit}P(y_i = 1) = \alpha_d + X_i' \beta^{(d)}, \quad d = 1, 2$$

- ▶ Second approach:

$$y = \dot{g}(\alpha_d + X\beta^{(d)}) + \varepsilon_d, \quad d = 1, 2$$

where $g(u) = \log(e^u + e^{-u})$ and ε_d is sub-Gaussian.

Part II:

CHIME: Clustering of High-dimensional Gaussian Mixtures with EM

Clustering

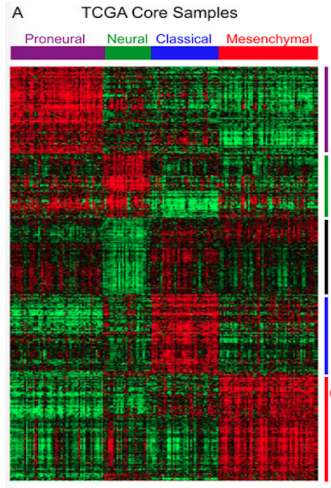


Fig: Verhaak et al. Cancer Cell, 2010

Clustering - existing algorithms

- ▶ K-means / K-median

Clustering - existing algorithms

- ▶ K-means / K-median
- ▶ Hierarchical clustering

Clustering - existing algorithms

- ▶ K-means / K-median
- ▶ Hierarchical clustering
- ▶ Expectation-Maximization (EM) algorithm

Clustering - existing algorithms

- ▶ K-means / K-median
- ▶ Hierarchical clustering
- ▶ Expectation-Maximization (EM) algorithm
- ▶ ...

Clustering - existing algorithms

- ▶ K-means / K-median
- ▶ Hierarchical clustering
- ▶ Expectation-Maximization (EM) algorithm
- ▶ ...

However, theoretical performance of the clustering algorithm is not fully understood.

Gaussian mixture model

- General form (2-class):

$$y^{(1)}, \dots, y^{(n)} \text{ i.i.d. } \sim \begin{cases} 1, & \text{with probability } 1 - \omega; \\ 2, & \text{with probability } \omega. \end{cases}$$

$$\mathbf{z}^{(i)} \mid y^{(i)} = d \text{ i.i.d. } \sim N_p(\boldsymbol{\mu}_d, \Sigma); \quad d = 1, 2.$$

Gaussian mixture model

- General form (2-class):

$$y^{(1)}, \dots, y^{(n)} \text{ i.i.d. } \sim \begin{cases} 1, & \text{with probability } 1 - \omega; \\ 2, & \text{with probability } \omega. \end{cases}$$

$$\mathbf{z}^{(i)} \mid y^{(i)} = d \text{ i.i.d. } \sim N_p(\boldsymbol{\mu}_d, \Sigma); \quad d = 1, 2.$$

- Observations: $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(n)}$.

Gaussian mixture model

- ▶ General form (2-class):

$$y^{(1)}, \dots, y^{(n)} \text{ i.i.d. } \sim \begin{cases} 1, & \text{with probability } 1 - \omega; \\ 2, & \text{with probability } \omega. \end{cases}$$

$$\mathbf{z}^{(i)} \mid y^{(i)} = d \text{ i.i.d. } \sim N_p(\boldsymbol{\mu}_d, \Sigma); \quad d = 1, 2.$$

- ▶ Observations: $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(n)}$.
- ▶ Goal: Cluster $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)}$ into two groups with **statistical guarantees**.

Gaussian mixture model

- ▶ When p is small, we solve for MLE to maximize

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \log \left\{ f(\mathbf{z}^{(i)} | \mu_1, \Sigma) P(y^{(i)} = 1) + f(\mathbf{z}^{(i)} | \mu_2, \Sigma) P(y^{(i)} = 2) \right\}.$$

⁴ Dempster et al. JRSSB. 1977

Gaussian mixture model

- ▶ When p is small, we solve for MLE to maximize

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \log \left\{ f(\mathbf{z}^{(i)} | \mu_1, \Sigma) P(y^{(i)} = 1) + f(\mathbf{z}^{(i)} | \mu_2, \Sigma) P(y^{(i)} = 2) \right\}.$$

- ▶ **Drawbacks:** $L(\theta)$ is not convex; MLE is challenging for large p .

⁴ Dempster et al. JRSSB. 1977

Gaussian mixture model

- ▶ When p is small, we solve for MLE to maximize

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \log \left\{ f(\mathbf{z}^{(i)} | \mu_1, \Sigma) P(y^{(i)} = 1) + f(\mathbf{z}^{(i)} | \mu_2, \Sigma) P(y^{(i)} = 2) \right\}.$$

- ▶ **Drawbacks:** $L(\theta)$ is not convex; MLE is challenging for large p .
- ▶ **Solution:** Expectation-Maximization (EM) algorithm⁴.

⁴ Dempster et al. JRSSB. 1977

Linear discriminant analysis

- Suppose we know the true parameters ω , μ_1 , μ_2 and Σ .

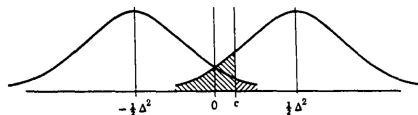


Fig: Mis-classification error of LDA⁵

⁵ $\Delta^2 = (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)$ is the SNR.

Linear discriminant analysis

- ▶ Suppose we know the true parameters ω , μ_1 , μ_2 and Σ .
- ▶ The discriminating direction $\beta = \Sigma^{-1}(\mu_1 - \mu_2)$.

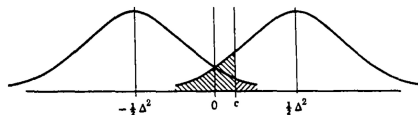


Fig: Mis-classification error of LDA⁵

⁵ $\Delta^2 = (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2)$ is the SNR.

Linear discriminant analysis

- ▶ Suppose we know the true parameters ω , μ_1 , μ_2 and Σ .
- ▶ The discriminating direction $\beta = \Sigma^{-1}(\mu_1 - \mu_2)$.
- ▶ Then the optimal classification rule

$$C_{opt}(\mathbf{z}) = \begin{cases} 1, & \{\mathbf{z} - (\mu_1 + \mu_2)/2\}'\beta \geq \log(\frac{\omega}{1-\omega}) \\ 2, & \{\mathbf{z} - (\mu_1 + \mu_2)/2\}'\beta < \log(\frac{\omega}{1-\omega}). \end{cases}$$

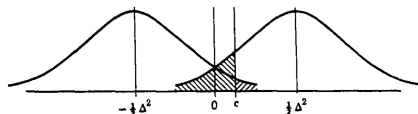


Fig: Mis-classification error of LDA⁵

⁵ $\Delta^2 = (\mu_1 - \mu_2)'\Sigma^{-1}(\mu_1 - \mu_2)$ is the SNR.

Linear programming discriminant

- If we know the sample labels $y^{(1)}, \dots, y^{(n)}$, estimate μ_d by

$$\hat{\mu}_d = \frac{1}{n_d} \sum_{i=1}^n \mathbf{z}^{(i)} I(y^{(i)} = d), \quad d = 1, 2,$$

and

$$\hat{\Sigma} = \frac{1}{n} (n_1 \hat{\Sigma}_1 + n_2 \hat{\Sigma}_2).$$

⁶ Cai and Liu. JASA. 2011

Linear programming discriminant

- ▶ If we know the sample labels $y^{(1)}, \dots, y^{(n)}$, estimate μ_d by

$$\hat{\mu}_d = \frac{1}{n_d} \sum_{i=1}^n \mathbf{z}^{(i)} I(y^{(i)} = d), \quad d = 1, 2,$$

and

$$\hat{\Sigma} = \frac{1}{n} (n_1 \hat{\Sigma}_1 + n_2 \hat{\Sigma}_2).$$

- ▶ Assuming sparse β , apply the LPD⁶ to get

$$\hat{\beta} = \arg \min \{ \|\beta\|_1 : \|\hat{\Sigma}\beta - (\hat{\mu}_1 - \hat{\mu}_2)\|_\infty \leq \lambda_n \}.$$

⁶ Cai and Liu. JASA. 2011

The EM algorithm

- ▶ We combine the above ideas to iteratively estimate $\theta = (\omega, \mu_1, \mu_2, \beta)$.

The EM algorithm

- ▶ We combine the above ideas to iteratively estimate $\theta = (\omega, \mu_1, \mu_2, \beta)$.
- ▶ The conditional log-likelihood

$$\begin{aligned} Q_n(\theta \mid \tilde{\theta}) &= \mathbb{E}_n[\log L(\theta; \tilde{\theta}, \mathbf{z})] \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{d=1}^2 P(y^{(i)} = d \mid \tilde{\theta}) \log f(\mathbf{z}^{(i)} \mid \mu_d, \Sigma) \end{aligned}$$

The EM algorithm

- ▶ Initialization $\theta^{(0)} = \{\omega^{(0)}, \boldsymbol{\mu}_d^{(0)}, \boldsymbol{\beta}^{(0)}\}; \kappa \in [1/2, 3/4]; \lambda_n^{(0)}; T_{stop}$.

The EM algorithm

- ▶ Initialization $\theta^{(0)} = \{\omega^{(0)}, \mu_d^{(0)}, \beta^{(0)}\}; \kappa \in [1/2, 3/4]; \lambda_n^{(0)}; T_{stop}$.
- ▶ E-step: Evaluate $Q_n(\theta \mid \theta^{(t)})$.

The EM algorithm

- ▶ Initialization $\theta^{(0)} = \{\omega^{(0)}, \mu_d^{(0)}, \beta^{(0)}\}; \kappa \in [1/2, 3/4]; \lambda_n^{(0)}; T_{stop}$.
- ▶ E-step: Evaluate $Q_n(\theta \mid \theta^{(t)})$.
- ▶ M-step:

$$(\omega^{(t+1)}, \mu_d^{(t+1)}, \Sigma^{(t+1)}) = \arg \max Q_n(\theta \mid \theta^{(t)}),$$
$$\lambda_n^{(t+1)} = \kappa \lambda_n^{(t)} + C \sqrt{\log p / n},$$

$$\beta^{(t+1)} = \arg \min \{ \|\beta\|_1 : \|\Sigma^{(t+1)}\beta - (\mu_1^{(t+1)} - \mu_2^{(t+1)})\|_\infty \leq \lambda_n^{(t+1)} \}.$$

The EM algorithm

- ▶ Initialization $\theta^{(0)} = \{\omega^{(0)}, \mu_d^{(0)}, \beta^{(0)}\}; \kappa \in [1/2, 3/4]; \lambda_n^{(0)}; T_{stop}$.
- ▶ E-step: Evaluate $Q_n(\theta \mid \theta^{(t)})$.
- ▶ M-step:

$$(\omega^{(t+1)}, \mu_d^{(t+1)}, \Sigma^{(t+1)}) = \arg \max Q_n(\theta \mid \theta^{(t)}),$$
$$\lambda_n^{(t+1)} = \kappa \lambda_n^{(t)} + C \sqrt{\log p / n},$$

$$\beta^{(t+1)} = \arg \min \{\|\beta\|_1 : \|\Sigma^{(t+1)}\beta - (\mu_1^{(t+1)} - \mu_2^{(t+1)})\|_\infty \leq \lambda_n^{(t+1)}\}.$$

- ▶ Upon convergence, output $\hat{\omega}, \hat{\mu}_d, \hat{\beta} \leftarrow \omega^{(T_{stop})}, \mu_d^{(T_{stop})}, \beta^{(T_{stop})}$.

Upper bound

Theorem (Cai, M, Zhang. 2018)

Assume $\|\beta\|_0 \leq s$. *Under certain technical conditions*, the output $\beta^{(T_{stop})}$ satisfies with high probability

$$\|\beta^{(T_{stop})} - \beta\|_2 \lesssim \kappa^{T_{stop}} \|\theta^{(0)} - \theta\|_2 + \sqrt{\frac{s \log p}{n}}.$$

Consequently, if $T_{stop} \gtrsim \log n$, then

$$\|\beta^{(T_{stop})} - \beta\|_2 \lesssim \sqrt{\frac{s \log p}{n}}.$$

Remarks

The results in Wang et al. ('15) ($\Sigma = \sigma^2 \mathbf{I}_p$) show

$$\|\beta^{(T_{stop})} - \beta\|_2 \lesssim \sqrt{\frac{s \log p \cdot \log n}{n}}.$$

The proposed classifier

- ▶ Given the estimated $\hat{\omega}$, $\hat{\mu}_d$, $\hat{\beta}$, the sample \mathbf{z} is classified as

$$\hat{C}(\mathbf{z}) = \begin{cases} 1, & \{\mathbf{z} - (\hat{\mu}_1 + \hat{\mu}_2)/2\}'\hat{\beta} \geq \log(\frac{\hat{\omega}}{1-\hat{\omega}}) \\ 2, & \{\mathbf{z} - (\hat{\mu}_1 + \hat{\mu}_2)/2\}'\hat{\beta} < \log(\frac{\hat{\omega}}{1-\hat{\omega}}). \end{cases}$$

The proposed classifier

- ▶ Given the estimated $\hat{\omega}, \hat{\mu}_d, \hat{\beta}$, the sample \mathbf{z} is classified as

$$\hat{C}(\mathbf{z}) = \begin{cases} 1, & \{\mathbf{z} - (\hat{\mu}_1 + \hat{\mu}_2)/2\}'\hat{\beta} \geq \log(\frac{\hat{\omega}}{1-\hat{\omega}}) \\ 2, & \{\mathbf{z} - (\hat{\mu}_1 + \hat{\mu}_2)/2\}'\hat{\beta} < \log(\frac{\hat{\omega}}{1-\hat{\omega}}). \end{cases}$$

- ▶ The mis-clustering error is defined as

$$R(\hat{C}) = \min_{\pi \in \mathbb{P}_2} \mathbb{E}[I(\hat{C}(\mathbf{z}) \neq \pi(y))],$$

where $\mathbb{P}_2 = \{\pi : [1, 2] \rightarrow [1, 2]\}$ is a set of permutation function.

Mis-clustering error

Theorem (Cai, M, Zhang. 2018)

Under the same conditions of Theorem 1 and with $T_{stop} \gtrsim \log n$, the classifier \hat{C} with mis-clustering error $R(\hat{C})$, satisfies w.h.p.

$$R(\hat{C}) - R_{opt} \lesssim \frac{s \log p}{n}.$$

Simulation

Competitors

- ▶ KM: k -means

Simulation

Competitors

- ▶ KM: k -means
- ▶ SKM: sparse k -means (Witten and Tibshirani '12)

Competitors

- ▶ KM: k -means
- ▶ SKM: sparse k -means (Witten and Tibshirani '12)
- ▶ SHP: sparse clustering with HARDT-PRICE (Azizyan et al. '14)

Competitors

- ▶ KM: k -means
- ▶ SKM: sparse k -means (Witten and Tibshirani '12)
- ▶ SHP: sparse clustering with HARDT-PRICE (Azizyan et al. '14)
- ▶ PCCM: penalized clustering with common covariances (Zhou et al. '09)

Simulation

Competitors

- ▶ KM: k -means
- ▶ SKM: sparse k -means (Witten and Tibshirani '12)
- ▶ SHP: sparse clustering with HARDT-PRICE (Azizyan et al. '14)
- ▶ PCCM: penalized clustering with common covariances (Zhou et al. '09)

Benchmark

- ▶ LPD: supervised linear program discriminant rule (Cai and Liu '11)

Simulation

Competitors

- ▶ KM: k -means
- ▶ SKM: sparse k -means (Witten and Tibshirani '12)
- ▶ SHP: sparse clustering with HARDT-PRICE (Azizyan et al. '14)
- ▶ PCCM: penalized clustering with common covariances (Zhou et al. '09)

Benchmark

- ▶ LPD: supervised linear program discriminant rule (Cai and Liu '11)
- ▶ Oracle: Fisher's LDA with true parameters

Simulation

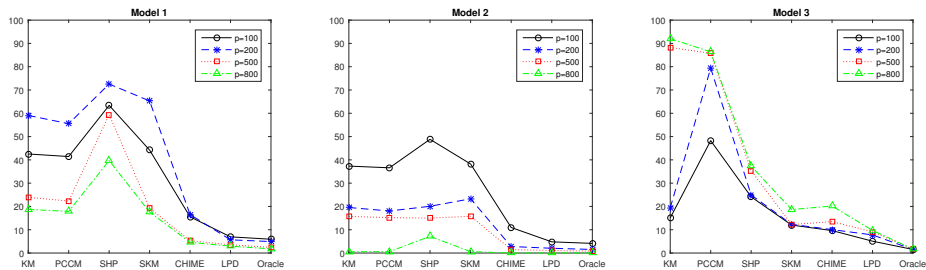


Fig: Clustering errors based on $n = 200$ test samples and 100 replications.

Application to clustering GBM data

Class	CHIME		KM		PCCM		SHP		SKM	
	1	2	1	2	1	2	1	2	1	2
Neural	26	0	26	0	26	0	12	14	25	1
Mesenchymal	2	54	7	49	5	51	10	46	6	50

Table: Clustering results for the GBM gene expression data with $p = 200$ genes and 82 samples

Summary



Summary

- ▶ Knowing labels doesn't improve the convergence rate of estimation and classification.



Summary



- ▶ Knowing labels doesn't improve the convergence rate of estimation and classification.
- ▶ Not covered in this talk
 - ▶ Lower bound of estimation and clustering error is in the same order of the respective upper bound.

Summary



- ▶ Knowing labels doesn't improve the convergence rate of estimation and classification.
- ▶ Not covered in this talk
 - ▶ Lower bound of estimation and clustering error is in the same order of the respective upper bound.
 - ▶ Extensions to multi-class GMM and/or unequal covariance matrices are available.

Collaborators

- ▶ Nanxun Ma (UW)
- ▶ Michael Wu (Hutch)
- ▶ Linjun Zhang (U Penn)
- ▶ Tony Cai (U Penn)