
Statistical Learning and Data Science

Xinwei Deng

xdeng@vt.edu

Department of Statistics

Course Agenda

- Presentation and discussion
- Choose topics and literature review
- Conduct a mini-research project

Statistical Perspectives for Modern Data Analytics

- Matrix: an arrangement of multiple rows (or columns)

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ a_{31} & a_{32} & \cdots & a_{3n} \\ \cdots & \cdots & & \cdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

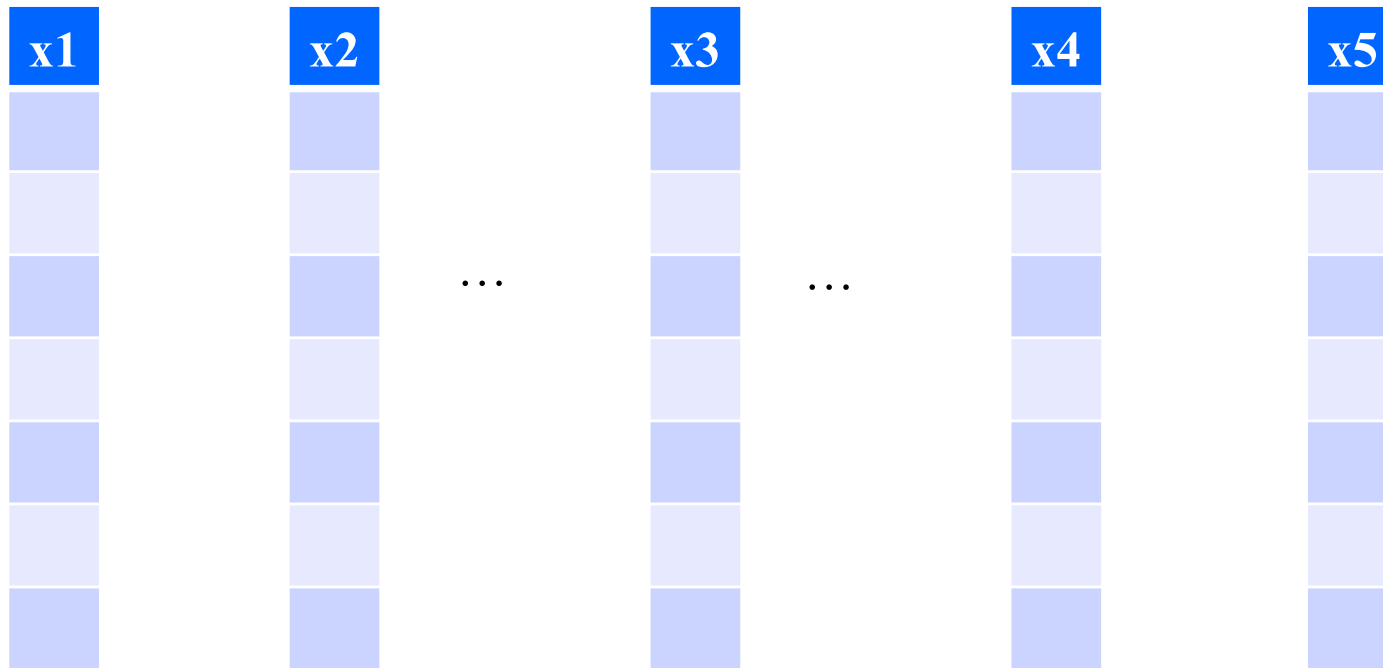
- From statistical perspective, it is mainly for **data matrix, or some statistic of data matrix.**
 - E.g., Covariance matrix
- The matrix algebra is useful and meaningful in modern data analytics.
 - E.g., matrix completion problem.

Statistical Perspectives for Modern Data Analytics (Con't)

- Big Data: Data is so-called Big with multiple rows, multiple columns, multiple types, multiple collection channels, etc.
- From statistical perspective, the challenges in Big Data mainly come from **dependency**.
 - Dependency among rows.
 - Dependency among columns.
 - Dependency between past-time and current-time data points.
 - Dependency between computation efficiency and estimation accuracy.

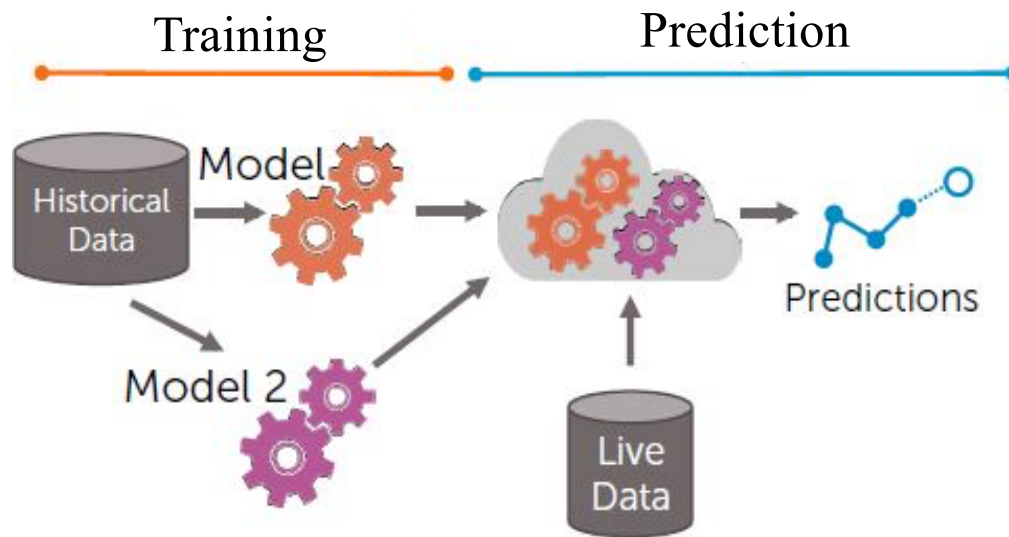
Illustration of Dependency

- In statistics, the challenges in Big Data mainly come from **dependency**.
 - Dependency among rows and dependency among columns.



Statistical Perspectives for Modern Data Analytics (Con't)

- Data Analytics: In simple view, it is similar to machine learning.



- From statistical perspective, the key of data analytics is to understand the **data generation mechanism**:
 - Enable model inference.
 - Enable model prediction.

Potential Topics 1

- Multi-response Regression and Model Selection
 - Structured Lasso with multiple responses.
 - Multi-response with consideration of covariance matrix.

Potential Topics 2

- Discriminant Analysis (DA) for Classification
 - LDA and QDA.
 - Regularized discriminant analysis.
 - Discriminant analysis based approach for classification.

Potential Topics 3

- Graphical Model via Covariance Matrix Estimation
 - Modified Cholesky decomposition approach.
 - Structured graphical model estimation.
 - Multiple graphical models.
 - Nonparametric approach for graphical model estimation.

Potential Topics 4

- Spatial-Temporal Modeling
 - Gaussian process modeling in computer experiment.
 - Multivariate time series with spatial correlation structure.

Potential Topics 5

- Interface between machine learning and experimental design
 - Effective data collection for efficient modeling.
 - Active learning and sequential design.
 - Embrace experimental design thinking for large-scale statistical analysis.

Potential Topics 6

- Statistical computation thinking in multivariate analysis
 - Tradeoff between computation efficiency and estimation accuracy.
 - Optimization-based data analytics.
 - Statistical computation thinking for data science.

Scope of Data and Modeling

- Scalar \rightarrow Vector \rightarrow Matrix \rightarrow Tensor
 - Many common concepts.
- From simple to advance
 - Simple regression to regularized tensor regression
 - Old techniques can be original.
- The responses can be continuous, discrete, or both.
 - From logistic regression to image classification.
- Modeling strategy:
 - Global vs Local modeling
 - Parametric vs Nonparametric methods.

The Multivariate Data

y_1	y_2	y_3	y_4	\dots	y_q

$\mathbf{Y} = (y_{ij})$ is $n \times q$ output matrix

x_1	x_2	x_3	\dots	x_k	x_{k+1}	\dots	\dots	x_{p-1}	x_p

$\mathbf{X} = (x_{ij})$ is $n \times p$ input matrix

Temporal/Functional Multivariate Data

y_1	y_2	y_3	y_4	...	y_q

$\mathbf{Y}(t) = (y_{ij})$ is $n \times q$
output matrix

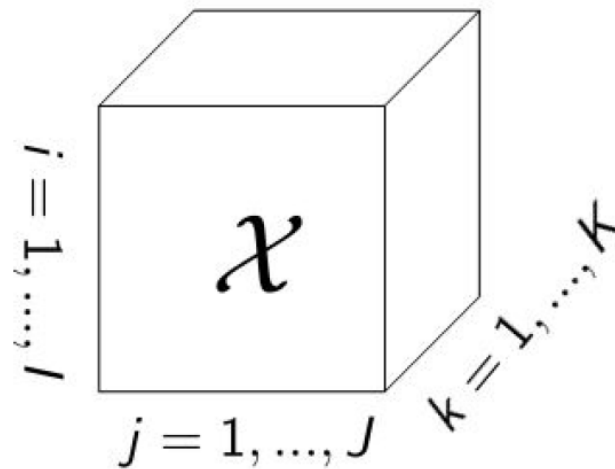
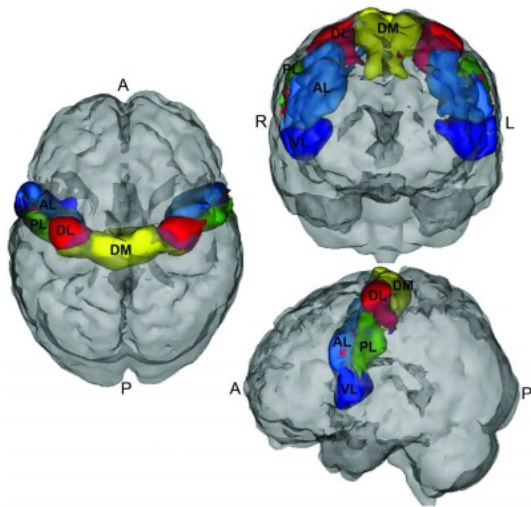
x_1	x_2	x_3	...	x_k	x_{k+1}	x_{p-1}	x_p

$\mathbf{X}(t) = (x_{ij})$ is $n \times p$ input
matrix

- At each time point, it is a multivariate data.

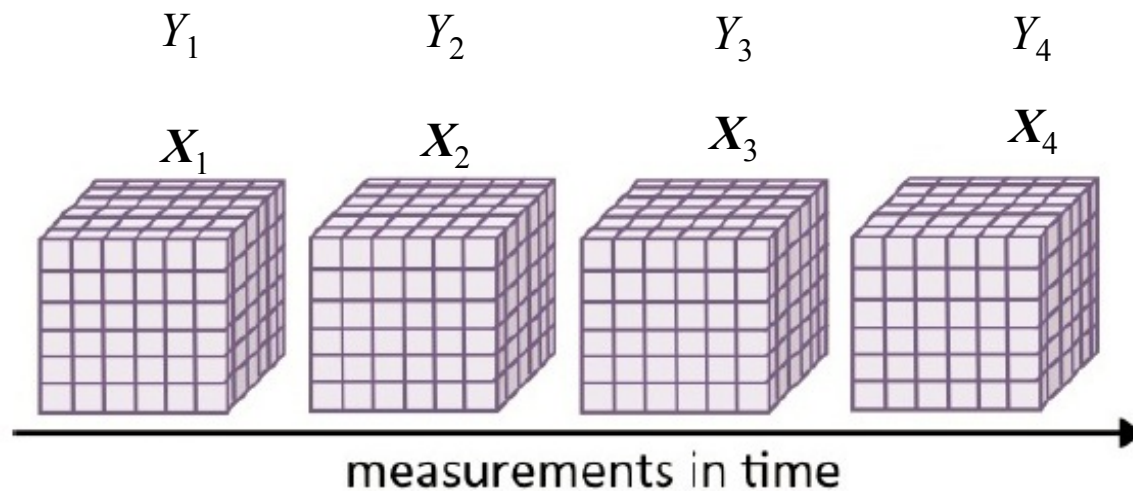
Tensor: a Generalization of Matrix

- Tensor: a *multidimensional* array.
- Example: fMRI data in Neuroscience: 3D brain images.
 - seek association between brain images (\mathbf{X}) and clinical outcomes (\mathbf{Y}).



Tensor Regression

- Tensor along time is a generalization of spatial-temporal structure.
 - Goal: predict the response on the future time point.



Topic 1:

Multi-response Regression

 y_1 \sim x_1 x_2 x_3 \dots x_k x_{k+1} \dots \dots x_{p-1} x_p \vdots y_m \sim x_1 x_2 x_3 \dots x_k x_{k+1} \dots \dots x_{p-1} x_p \vdots y_q \sim x_1 x_2 x_3 \dots x_k x_{k+1} \dots \dots x_{p-1} x_p

The Linear Model

- ▶ Consider the linear model $y = x'\beta + \epsilon$, $\epsilon \sim N(0, \sigma^2)$ where $x = (x_1, \dots, x_p)'$ and $\beta = (\beta_1, \dots, \beta_p)'$.
- ▶ With data $(x_i, y_i), i = 1, \dots, n$, the log-likelihood is

$$\begin{aligned} L(\beta, \sigma^2) &= \log \left\{ \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} (1/\sigma^2)^{1/2} \exp\left(-\frac{(y_i - x_i'\beta)^2}{2\sigma^2}\right) \right\} \\ &\propto -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - x_i'\beta)^2 \\ &= -\frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta) \end{aligned}$$

- ▶ The MLE of β is obtained by $\min -\log(L(\beta, \sigma^2))$,

$$\min_{\beta, \sigma^2} \frac{n}{2} \log \sigma^2 + \frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta).$$

Some Remarks

- ▶ Estimating β in MLE is equivalent to the OLS estimation, i.e.,

$$\min_{\beta} LS(\beta) = (y - X\beta)^T (y - X\beta)$$

- ▶ The negative log-likelihood or least squares both can be viewed as loss functions.
- ▶ Consider the regularization/penalty on the loss function, which is in the format of

$$\text{Loss Function} + \lambda \text{Penalty}$$

- ▶ The penalized likelihood approach is the same as the penalized least squares when the penalty only involves β .

Lasso for Model Selection

- ▶ The Lasso is to estimate β by

$$\hat{\beta}^{\text{Lasso}} = \operatorname{argmin}_{\beta} \left\{ (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \sum_{j=1}^p |\beta_j| \right\}.$$

- ▶ Here is the penalty $P(\beta)$ is l_1 norm of β , i.e.
 $P(\beta) = \|\beta\|_1 = \sum_{j=1}^p |\beta_j|.$
- ▶ It is equivalent to the constraint/regularization problem:

$$\begin{aligned} \hat{\beta}^{\text{Lasso}} &= \operatorname{argmin}_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \\ &\text{subject to } \sum_{j=1}^p |\beta_j| \leq M \end{aligned}$$

- ▶ The other related: adaptive Lasso, fused Lasso, and the generalized Lasso.

Group Lasso for Model Selection

- ▶ Consider the predictor variables have a group structure by

$$\mathbf{x} = (x_{11}, \dots, x_{1k_1}, \dots, x_{p1}, \dots, x_{pk_p}),$$

where each group is x_{j1}, \dots, x_{jk_j} .

- ▶ Such group structures can be non-overlapped or overlapped.
- ▶ The linear model can be written as

$$y = \sum_{i=1}^p \sum_{j=1}^{k_i} x_{ij} \beta_{ij} + \epsilon.$$

- ▶ The penalty term also accommodate the group structure:
 - ▶ Yuan and Lin (2006): $P(\boldsymbol{\beta}) = \sum_{j=1}^p \sqrt{\beta_{j1}^2 + \dots + \beta_{jk_j}^2}$.
 - ▶ Zhao et al. (2006): $P(\boldsymbol{\beta}) = \sum_{j=1}^p \max\{|\beta_{j1}|, \dots, |\beta_{jk_j}|\}$.
 - ▶ Group Bridge (Huang et. al., 2009) and Hierarchical LASSO (Zhou and Zhu, 2009).

Multi-response Linear Model

- ▶ In multi-response regression, the response vector $y = (y_1, \dots, y_q)'$ and the predictor vector $x = (x_1, \dots, x_p)'$.
- ▶ Consider the linear model

$$y = B'x + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I),$$

where B is a $p \times q$ matrix of coefficients and the k th column β_k is coefficient vector associated with y_k regressing on x .

- ▶ Remark 1: For these q regression models, different response variable with the same predictor variables.
- ▶ Remark 2: One can also consider the model $y|x \sim N(B'x, \Sigma)$. Here Σ reflects the correlation structure among q response variables $y = (y_1, \dots, y_q)'$.

Model Selection for Multi-response Model

- ▶ The log-likelihood function $L(B)$ is

$$\log \left\{ \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} |\Sigma^{-1}|^{1/2} \exp\left(-\frac{(y_i - B'x_i)' \Sigma^{-1} (y_i - B'x_i)}{2}\right) \right\}$$
$$\propto \frac{n}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n (y_i - B'x_i)' \Sigma^{-1} (y_i - B'x_i)$$

with $\Sigma = \sigma^2 I$.

- ▶ Seeking sparse estimation of coefficient matrix B , consider penalized regression based on the log-likelihood function.

$$\hat{B} = \arg \min \frac{1}{n} \sum_{i=1}^n (y_i - B'x_i)' (y_i - B'x_i) + \lambda P(B)$$
$$= \arg \min \sum_{k=1}^q \frac{1}{n} \sum_{i=1}^n (y_{ik} - \beta'_k x_i)^2 + \lambda P(B).$$

Penalty: Structured v.s. Non-Structured

- ▶ The coefficient matrix B is in a matrix format as

$$B = \begin{pmatrix} \beta_{11} & \dots & \beta_{q1} \\ \vdots & \ddots & \vdots \\ \beta_{1p} & \dots & \beta_{qp} \end{pmatrix}.$$

- ▶ One can pursue
 - ▶ A sparse coefficient matrix estimate \hat{B} .
 - ▶ Some rows of \hat{B} become zero vectors.
 - ▶ Other patterns of interest on the sparsity.
- ▶ To encourage some rows of B being zeros for reducing number of predictors in the model, consider the penalty

$$P(B) = \lambda_1 \sum_{j=1}^p \sqrt{\sum_{k=1}^q \beta_{kj}^2} + \lambda_2 \sum_{j=1}^p \sum_{k=1}^q |\beta_{kj}|.$$

where λ_1 and λ_2 are tuning parameters.

Penalty: Hierarchical Structure via Re-parameterization

- ▶ Alternatively, we can consider a hierarchical structure to parameterize β_{kj} as

$$\beta_{jk} = \gamma_j \alpha_{kj}, \quad \gamma_j \geq 0.$$

- ▶ Such a parametrization provides the flexibility of obtaining a penalty with pursuing sparsity along the rows of B .
- ▶ Specifically, we can consider the penalty function as follows:

$$P(B) = \lambda_1 \sum_{j=1}^p \gamma_j + \lambda_2 \sum_{j=1}^p \sum_{k=1}^q |\alpha_{kj}|,$$

where λ_1 and λ_2 are tuning parameters.

Multi-response with Covariance Matrix Estimation

- ▶ The multi-response linear model $\mathbf{y} = \mathbf{B}'\mathbf{x} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \boldsymbol{\Sigma})$.
- ▶ To estimate \mathbf{B} and $\boldsymbol{\Sigma}$, consider the negative log-likelihood which is

$$\begin{aligned} & -\frac{n}{2} \log |\boldsymbol{\Sigma}^{-1}| + \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \mathbf{B}'\mathbf{x}_i)' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \mathbf{B}'\mathbf{x}_i) \\ & = -\frac{n}{2} \log |\boldsymbol{\Sigma}^{-1}| + \frac{1}{2} \text{tr}[(\mathbf{Y} - \mathbf{X}\mathbf{B})\boldsymbol{\Sigma}^{-1}(\mathbf{Y} - \mathbf{X}\mathbf{B})^T] \end{aligned}$$

- ▶ Seeking the sparsity on \mathbf{B} with the consideration of $\boldsymbol{\Sigma}$, we can consider

$$\begin{aligned} \min & -\log |\boldsymbol{\Sigma}^{-1}| + \frac{1}{n} \text{tr}[(\mathbf{Y} - \mathbf{X}\mathbf{B})\boldsymbol{\Sigma}^{-1}(\mathbf{Y} - \mathbf{X}\mathbf{B})^T] \\ & + \lambda_1 P_1(\mathbf{B}) + \lambda_2 P_2(\boldsymbol{\Sigma}) \end{aligned}$$

A few Remarks

- The complex of penalty is closely related to the complex of the model estimation.
- The MRCE (Rothman et al., 2010) is to consider an overall sparse structure on \mathbf{B} and Σ^{-1} , which is

$$P_1(\mathbf{B}) = \sum_{i,j} |b_{ij}|, P_2(\Sigma^{-1}) \sum_{i \neq j} |c_{ij}|,$$

where c_{ij} is the (i, j) element of Σ^{-1} .

- Therefore, the MRCE is to consider

$$\begin{aligned} \min & -\log |\Sigma^{-1}| + \frac{1}{n} \text{tr}[(\mathbf{Y} - \mathbf{X}\mathbf{B})\Sigma^{-1}(\mathbf{Y} - \mathbf{X}\mathbf{B})^T] \\ & + \lambda_1 \sum_{i,j} |b_{ij}| + \lambda_2 \sum_{i \neq j} |c_{ij}|. \end{aligned}$$

Another Angle from Multivariate t-Distribution

- ▶ The multivariate t-distribution can be viewed as a scaled multivariate normal.
- ▶ Specifically, denote $w \sim \Gamma(\nu/2, \nu/2)$, which is independent of $z \sim N(\mathbf{0}, \Sigma)$. If defining

$$\mathbf{y} = \boldsymbol{\mu} + w^{-1/2} \mathbf{z},$$

Then $y \sim t(\boldsymbol{\mu}, \Sigma, \nu)$.

- ▶ Using the above result, one can develop multi-response regression using t-distribution as the random error, i.e.,

$$\mathbf{y} = \mathbf{B}'\mathbf{x} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim t(\mathbf{0}, \Sigma, \nu).$$

Another Angle from Multivariate t-Distribution (Con't)

- ▶ Under t-distribution as the random error, the multi-response linear model is

$$\mathbf{y} \sim t(\mathbf{B}'\mathbf{x}, \mathbf{\Sigma}, \nu) \Leftrightarrow \sqrt{w}(\mathbf{y} - \mathbf{B}'\mathbf{x})|w \sim N(0, \mathbf{\Sigma})$$

- ▶ In addition, one can show that $w|\mathbf{y}$ still has a Gamma distribution, which is

$$\begin{aligned} w|\mathbf{y} &\sim \Gamma\left(\frac{\nu + p}{2}, \frac{\nu + d_{\mathbf{y}}(\boldsymbol{\mu}, \mathbf{\Sigma})}{2}\right) \\ \Rightarrow E(w|\mathbf{y}) &= \frac{\nu + p}{\nu + d_{\mathbf{y}}(\boldsymbol{\mu}, \mathbf{\Sigma})} \end{aligned}$$

where $d_{\mathbf{y}}(\boldsymbol{\mu}, \mathbf{\Sigma})$ is the square of Mahalanobis distance,
 $d_{\mathbf{y}}(\boldsymbol{\mu}, \mathbf{\Sigma}) = (\mathbf{y} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu})$.

- ▶ Therefore, we can develop an EM algorithm to obtain the parameter estimation.

EM Algorithm for Parameter Estimation

- ▶ Suppose the data are $(\mathbf{x}_i, \mathbf{y}_i), i = 1, \dots, n$, and $\mathbf{x} \in \mathcal{R}^p$, $\mathbf{y} \in \mathcal{R}^q$.
- ▶ With an initial estimation of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, we can develop an EM algorithm for parameter estimation. At k th iteration

E-STEP Given current estimator $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Sigma}}$, update the value of w_i ,

$$w_i^{(k)} = \frac{\nu + p}{\nu + d_i(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}})},$$

where $d_i(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) = (\mathbf{y} - \hat{\boldsymbol{\mu}})^T \hat{\boldsymbol{\Sigma}}^{-1} (\mathbf{y} - \hat{\boldsymbol{\mu}})$.

M-STEP With the updated value of $w_i^{(k)}$, we can update the estimation of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ by

$$\begin{aligned} \min & -\log |\boldsymbol{\Sigma}^{-1}| + \frac{1}{n} \sum_{i=1}^n w_i^{(k)} (\mathbf{y}_i - \mathbf{B}' \mathbf{x}_i)' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \mathbf{B}' \mathbf{x}_i) \\ & + \lambda_1 \sum_{i,j} |b_{ij}| + \lambda_2 \sum_{i \neq j} |c_{ij}|. \end{aligned}$$

Topics 3:

Covariance Matrix Estimation

- ▶ The multi-response model is simplified as

$$\mathbf{y} = \boldsymbol{\mu} + \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim N(0, \boldsymbol{\Sigma}).$$

- ▶ With data $\mathbf{y}_i, i = 1, \dots, n$, it is equivalent to

$$\mathbf{y}_i \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}), i = 1, \dots, n.$$

- ▶ To estimate $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$, the log-likelihood function becomes

$$\begin{aligned} L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \log \left\{ \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} |\boldsymbol{\Sigma}^{-1}|^{1/2} \exp\left(-\frac{(\mathbf{y}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu})}{2}\right) \right\} \\ &\propto \frac{n}{2} \log |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}) \\ &\propto \log |\boldsymbol{\Sigma}^{-1}| - \text{tr}[\boldsymbol{\Sigma}^{-1} \mathbf{S}], \end{aligned}$$

where

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})^T.$$

Estimating Σ

- Based on

$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \frac{n}{2} \log |\boldsymbol{\Sigma}^{-1}| - \frac{1}{2} \sum_{i=1}^n (\mathbf{y}_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_i - \boldsymbol{\mu})$, it is easy to obtain the estimate of $\boldsymbol{\mu}$ as

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i.$$

- To estimate $\boldsymbol{\Sigma}$, take the first derivative with respect to $\boldsymbol{\Omega} = \boldsymbol{\Sigma}^{-1}$

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\Omega}} \{ -\log |\boldsymbol{\Omega}| + \text{tr}[\boldsymbol{\Omega} \mathbf{S}] \} \\ = -\boldsymbol{\Sigma} + \mathbf{S} = 0 \end{aligned}$$

Therefore, the estimate of $\boldsymbol{\Sigma}$ is

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^T.$$

- What if taking the first derivative w.r.t. $\boldsymbol{\Sigma}$? (Take-home)

Why Estimating Σ^{-1} : Gaussian Graphical Model

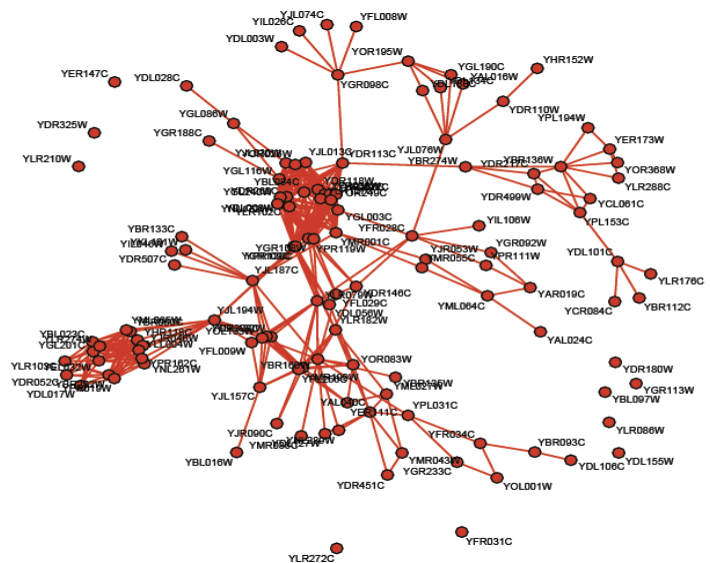
[illegible]

$$\mathbf{\Omega} = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1q} \\ c_{21} & c_{22} & \dots & c_{2q} \\ \vdots & \ddots & \vdots & \vdots \\ c_{q1} & c_{q2} & \dots & c_{qq} \end{pmatrix}$$

concentration matrix $\Sigma^{-1} \equiv \Omega = (c_{ij})$

Data $\mathbf{Y} = (y_{ij})$ is an $n \times q$ matrix

- describe the **conditional dependency** among variables: if $c_{ij}=0$ zero, then variables i and j are conditionally independent given the other variables.



Thank you!

- **Questions and Comments?**