

Final exam: Tuesday April 19 7-9 pm

Coverage: Ch 1-6 of text + classification (mainly LDA, QDA)

Allowed 2 8.5 x 11 sheet (both sides) & a calculator

Office hour: Next week Wed 10am - noon

Thurs 1-3 pm

Mon April 18 1-3 pm

Structured Multivariate Models

Last time: Repeated Measures

$$\underline{X}_i = \begin{pmatrix} X_{i1} \\ \vdots \\ X_{ik} \end{pmatrix} = \begin{pmatrix} \mu_i \\ \vdots \\ \mu_k \end{pmatrix} + A_i \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} + \begin{pmatrix} \varepsilon_{i1} \\ \vdots \\ \varepsilon_{ik} \end{pmatrix}$$

$\text{Var}(A_i) = \sigma_A^2$ $\text{Var}(\varepsilon_{ij}) = \sigma_\varepsilon^2$

Longitudinal data

n subjects measured at time t_1, \dots, t_k

$$X_{ij} = X_i(t_j) \Rightarrow \underline{X}_i = \begin{pmatrix} X_i(t_1) \\ \vdots \\ X_i(t_k) \end{pmatrix}$$

Subject

Examples

① Growth curves: normal vs abnormal growth

② Biological measurements: cholesterol level over time

Question: How do longitudinal data differ from, say, repeated measure data?

For each subject, we would expect that

$$\text{Var}(X_i(t_j) - X_i(t_k)) = \text{Var}(X_i(t_j)) + \text{Var}(X_i(t_k)) - 2\text{Cov}(X_i(t_j), X_i(t_k))$$

increases as $|t_j - t_k|$ increases

Equivalently, $\text{Cov}(X_i(t_j), X_i(t_k))$ decreases as $|t_j - t_k|$ increases

- need to incorporate this into modeling.

Modeling: Starting pt. $\xrightarrow{\text{random}} A_i$ $\xrightarrow{\text{fixed}} \mu_i(t_j)$ $\xrightarrow{\text{noise}} W_i(t_j)$

$$X_i(t_j) = A_i + \mu_i(t_j) + W_i(t_j)$$

or $X_i(t_j) = A_i + \mu_i(t_j) + W_i(t_j) + \varepsilon_{ij}$ \leftarrow measurement error

Assumptions:

① A_1, \dots, A_n are independent $N(0, \sigma_A^2)$

② W_1, \dots, W_n are independent stochastic process with $E[W_i(t)] = 0$
and $\text{Cov}(W_i(t_j), W_i(t_k)) = \gamma(|t_j - t_k|)$

$$\gamma(s) \rightarrow 0 \text{ as } s \rightarrow 0$$

\downarrow
autocovariance function

③ $\{\varepsilon_{ij}\}$ are independent $N(0, \sigma_\varepsilon^2)$

④ All random components are independent

$$\underline{X}_i = \begin{pmatrix} X_i(t_1) \\ \vdots \\ X_i(t_k) \end{pmatrix} \sim N_k(\mu_i, C)$$

where elements of μ_i, C are given by:

$$\mu = \begin{pmatrix} \mu(t_1) \\ \vdots \\ \mu(t_k) \end{pmatrix} \quad \text{Var}(X_i(t_j)) = E[(X_i(t_j) - \mu(t_j))^2]$$

$$\begin{aligned} \text{Var}(X_i(t_j)) &= \text{Var}(A_i + W_i(t_j) + \varepsilon_{ij}) \\ &= \text{Var}(A_i) + \text{Var}(W_i(t_j)) + \text{Var}(\varepsilon_{ij}) \\ &= \sigma_A^2 + \gamma(0) + \sigma_\varepsilon^2 \end{aligned}$$

$$X_i(t_j) = A_i + \mu(t_j) + W_i(t_j) + \varepsilon_{ij}$$

$$\begin{aligned} \text{Cov}(X_i(t_j), X_i(t_k)) &= E[(A_i + W_i(t_j) + \varepsilon_{ij})(A_i + W_i(t_k) + \varepsilon_{ik})] \\ &= \sigma_A^2 + \gamma(|t_j - t_k|) \end{aligned}$$

Need to specify $\mu \rightarrow \gamma(s)$ more precisely

Simple estimate for μ

$$\hat{\mu} = \begin{pmatrix} \hat{\mu}(t_1) \\ \vdots \\ \hat{\mu}(t_k) \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \tilde{X}_i \rightarrow \hat{\mu}(t_j) = \frac{1}{n} \sum_{i=1}^n X_i(t_j)$$

We can also specify $\mu(t)$ parametrically

$$\mu(t) = \beta_0 + \sum_{i=1}^n \beta_i \phi_i(t)$$

known
unknown

- ϕ_1, \dots, ϕ_p are some specified functions e.g. polynomials, sines, sinusoids.

Autocovariance function

- typically $\gamma(s) \rightarrow 0$ as $s \rightarrow \infty$
 - parametric family: $\gamma(s) = \gamma(0) \exp(-\alpha s^\lambda)$ $(\alpha > 0, \lambda > 0)$
- Unknown

- special case: $\lambda = 1 \Rightarrow$ autoregressive autocovariance function

One approach: Throw everything into likelihood function and maximize \Rightarrow black box

$$L(\sigma_A^2, \sigma_\varepsilon^2, \gamma(0), \alpha, \lambda, \beta_0, \beta_1, \dots, \beta_p)$$

Sample variogram: Estimating $\gamma(s)$

$$\text{What is it? Look at } \frac{1}{2} \text{Var}(X_i(t_j) - X_i(t_l)) = \underbrace{\frac{1}{2} \text{Var}(X_i(t_j)) + \frac{1}{2} \text{Var}(X_i(t_l))}_{-\text{Cov}(X_i(t_j), X_i(t_l))} \underbrace{\sigma_A^2 + \gamma(0) + \sigma_\varepsilon^2}_{\sigma_A^2 + \gamma(|t_j - t_l|)}$$

$$= \sigma_\varepsilon^2 + \gamma(0) - \gamma(|t_j - t_l|)$$

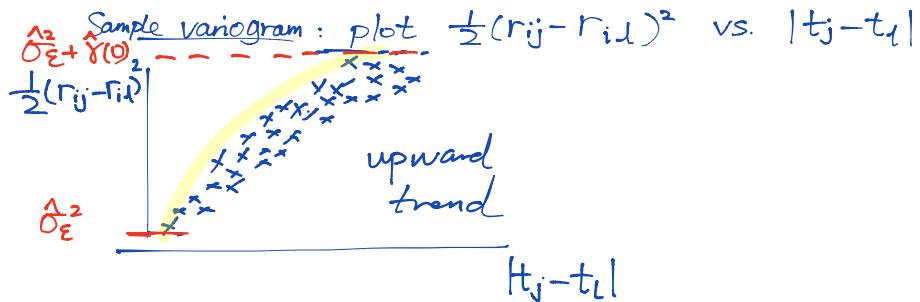
$$\begin{cases} = \sigma_\varepsilon^2 & \text{if } t_j = t_l \\ \approx \sigma_\varepsilon^2 + \gamma(0) & \text{if } |t_j - t_l| \text{ is sufficiently large} \end{cases}$$

前面的

$t_j, t_l \rightarrow$ 换成
 t_l

Given some estimate $\hat{\mu} = \begin{pmatrix} \hat{\mu}(t_1) \\ \vdots \\ \hat{\mu}(t_k) \end{pmatrix}$ define

$$\Gamma_{ij} = r_i(t_j) = x_i(t_j) - \hat{\mu}(t_j)$$



How to determine $\frac{\gamma(s)}{\gamma(0)} = \varphi(s)$?

$$\text{Spps } \gamma(s) = \gamma(0) \exp(-\alpha s^\lambda)$$

$$\therefore \varphi(s) = \exp(-\alpha s^\lambda)$$

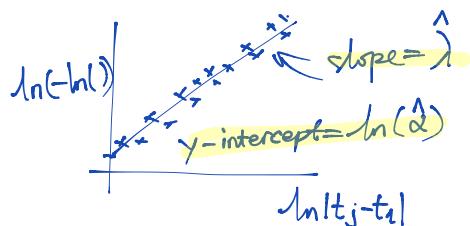
Estimate $\varphi(s)$ (i.e. α, λ) from variogram plot

$$\frac{\hat{\sigma}_\varepsilon^2 + \hat{\gamma}(0) - \frac{1}{2}(r_{ij} - r_{il})^2}{\hat{\gamma}(0)} \doteq \exp(-\alpha |t_{ij} - t_{il}|^\lambda)$$

$\approx \hat{\gamma}(|t_{ij} - t_{il}|)$

$$\ln \left(\frac{\hat{\sigma}_\varepsilon^2 + \hat{\gamma}(0) - \frac{1}{2}(r_{ij} - r_{il})^2}{\hat{\gamma}(0)} \right) \doteq -\alpha |t_{ij} - t_{il}|^\lambda$$

$$\ln(-\ln(\dots)) \doteq \ln(\alpha) + \lambda \ln |t_{ij} - t_{il}|$$



- use estimates to get std error for estimates of $\hat{\mu}(t_j)$

- use estimates from variogram as initial estimates in maximum likelihood estimation.