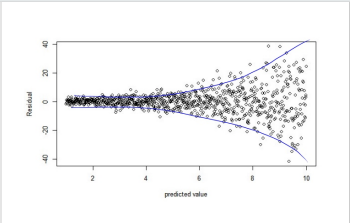
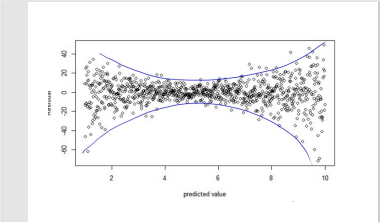


$$\sigma^2 g(\cdot) = \sigma^2 (\gamma_0 + \gamma_1 \hat{y})$$



$$\begin{aligned}\sigma^2 g(\cdot) &= \sigma^2 [e^{\gamma \hat{y}}] \\ \sigma^2 g(\cdot) &= \sigma^2 [\gamma_0 + e^{\gamma \hat{y}}]\end{aligned}$$



$$\sigma^2 g(\cdot) = \sigma^2 [\gamma_0 + \gamma_1 \hat{y} + \gamma_2 \hat{y}^2]$$

$\hat{y}_i = f(x_i, \theta)$ 이므로 결국 $g(\cdot)$ 도 θ 에 dependent 하다.

2. Statistical Modelling (2)

Statistical Modelling & Machine Learning

Jaejik Kim

Department of Statistics
Sungkyunkwan University

parametric 한 모델링에 중점

$$y = f(x) + \varepsilon$$

↑
Specified

STA3036

Data-Modelling Culture : parametric (small sample size, with information about the relationship between x and y)

Algorithmic Modelling Culture : non-parametric

Statistical Models

요즘 statistical modelling 이라고 한다면 Data-modelling 과 Algorithmic-modelling 을 합친 경우가 많으나, 수렴시문명 parametric 한 경우만 알아본다.

- ▶ Statistical (data) model: A method to look at data / Summary (reduction) of data. $E(X) \quad Y = \beta_0 + \beta_1 X$
n개의 데이터들 하나의 수식으로 표현한 것
- ▶ Statistical models consist of two elements; systematic and random effects. $E(X) \quad \beta_i$'s
 - ▶ Systematic effects: Pattern of data. $\beta_0, \beta_1, \beta_2, \dots$
 - ▶ Random effects: Unexplained or random variation. ϵ
 - ▶ Systematic effects are likely to be blurred by random effects.
 - ▶ Random effects are usually described in statistical terms.
확률적 분포
- ▶ Looking intelligently at data \Rightarrow Formulation of patterns \Rightarrow Statistical data models.
parametric
 - ▶ Succinct description of the systematic variation in the data.
 - ▶ Description patterns in similar data that might be collected for another study.
 \Rightarrow 분산 / 분포에 대한 자세한 설명
 \Rightarrow 유사한 데이터에 대한 설명력

Statistical Data Modelling (Parametric Models)

- ▶ E.g., consider the following model: $y = f(x; \theta)$.
 - ▶ No error & specified form of f .
 - ▶ For given x_1, \dots, x_n , y takes the values $f(x_1; \theta), \dots, f(x_n; \theta)$.
 - ▶ If θ is given, the values of y can be exactly reconstructed.
- ⇒ For given x_1, \dots, x_n , θ is an exact summary of y_1, \dots, y_n .

그러나 실제로

- ▶ Since there are errors in practice, the relationship between y and x has approximately f .
 - ▶ $\hat{y}_i = f(x_i; \hat{\theta})$, $i = 1, \dots, n$: Theoretical or fitted values generated by the model f and the data.
 - ▶ The model cannot reproduce the original data values y_1, \dots, y_n exactly. ↳ 가 확률변수를 따르기 때문에
 - ▶ The pattern from the model approximates the data values and it can be summarized by θ .

θ 를 어떻게 추정할 것인가

▶ Estimation methods for data models:

⇒ $n \rightarrow \infty$ 일 때, MLE는 consistent 하고, minimum variance를 가지며, normality를 가진다.

- ▶ **Maximum likelihood estimation:** Find model parameters maximizing the likelihood function for given data.

⇒ θ 가 fixed가 아닌 random variable로써 사용된다. (model-update에 용이하다) (불확실성을 수량화할 수 있다)

- ▶ **Bayesian estimation:** Find the posterior distribution of model parameters for given prior distributions and likelihood function.

⇒ 관찰되지 못한 불확실성에 대해 고려할 수 있다.

- ▶ This class focuses on the ML estimation.

Least Square Method

주된 machine learning 들의 가정: $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$

- ▶ Model: $Y = f(X; \theta) + \epsilon$.
 - ▶ Y : Continuous variable.
 - ▶ f : Model.
 - ▶ $X = (X_1, \dots, X_p)^\top$: Input variable vector.
 - ▶ θ Model parameter vector.
 - ▶ ϵ : Random error.
- ▶ Least square method: Find θ minimizing the discrepancy between y and \hat{y} .

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2,$$

where $\hat{y}_i = f(x_i; \hat{\theta})$.

- ▶ If (1) y_i 's are statistically independent and (2) the variance of y_i does not depend on its mean value, the LS criterion is valid as a measure of discrepancy between y and \hat{y} .
- ▶ Conditions (1) and (2) guarantee that all observations have the same weight. \Rightarrow Universal Variance를 각고, independent 하기 때문에

Maximum Likelihood Estimation

- ▶ Data: $(y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n)$.
- ▶ Assumption: X_1, \dots, X_p are given (constant).
- ▶ Regression function: $E(Y|X = x) = f(x; \theta)$.
constant value
- ▶ Random variables in the data: Y_1, \dots, Y_n .
- ▶ To construct the likelihood function, the joint distribution of Y_1, \dots, Y_n , $p(\mathbf{Y}; \theta)$, should be identified.
- ▶ Likelihood function:

$$L(\theta; \mathbf{y}) \equiv p(\mathbf{Y}; \theta).$$

joint distribution of Y_1, \dots, Y_n

- ▶ MLE of model parameter θ : Let $l(\theta) = \log L(\theta)$.
 - ▶ θ maximizing $l(\theta)$ or θ minimizing $-2l(\theta)$.
deviance

Relationship between LS and ML

- ▶ $\epsilon_i \sim^{iid} N(0, \sigma^2)$, $i = 1, \dots, n$.
- ▶ $Y_i \sim^{iid} N(\mu_i, \sigma^2)$, $i = 1, \dots, n$, where $\mu_i = E(Y_i) = f(\mathbf{x}_i; \boldsymbol{\theta})$.
- ▶ Since Y_i 's are independent, the joint density of $\mathbf{Y} = (Y_1, \dots, Y_n)^\top$ is

$\mu_i = f(\mathbf{x}_i; \boldsymbol{\theta})$ 가 주어졌으므로

$$\begin{aligned} p(\mathbf{Y}; \boldsymbol{\theta}) &= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(Y_i - \mu_i)^2}{2\sigma^2} \right) \right] \\ &= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(Y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2}{2\sigma^2} \right) \right]. \end{aligned}$$

Relationship between LS and ML

- MLE: For fixed σ^2 ,

$$\begin{aligned}\max_{\boldsymbol{\theta}} [l(\boldsymbol{\theta})] &\equiv \min_{\boldsymbol{\theta}} [-2l(\boldsymbol{\theta}; \mathbf{y})] \\ &\equiv \min_{\boldsymbol{\theta}} \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2 \\ &\equiv \min_{\boldsymbol{\theta}} \sum_{i=1}^n (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2 \\ &\equiv \min_{\boldsymbol{\theta}} (\mathbf{y} - \mathbf{f})^\top (\mathbf{y} - \mathbf{f}) \\ &= LS \text{ criterion},\end{aligned}$$

where $\mathbf{f} = (f(\mathbf{x}_1; \boldsymbol{\theta}), \dots, f(\mathbf{x}_n; \boldsymbol{\theta}))^\top$.

$\therefore \varepsilon_i \sim N(0, \sigma^2)$ 이 충족되면 MLE = LSE

When Error Assumptions are Violated

※ 일반적인 선형 회귀의 가정들이 위배되었을 때 어떻게 모델링을 해야 하나?

- Objective function을 Error Assumption이 위배되었을 때의 model로 modify해서 새로운 model을 유도.

▶ Assumptions for error ϵ :

- (1) ϵ_i 's have constant variance. Homoscedasticity
- (2) ϵ_i 's are independent.
- (3) ϵ_i 's have normal distribution.

residual plot을 통해
간단히 확인 가능

- ▶ From the residuals $r_i = y_i - \hat{y}_i$, $i = 1, \dots, n$, we can check the assumptions (1), (2) and (3).

- ▶ When the assumptions are violated, the model variance $\uparrow \Rightarrow$ Poor prediction.
※ n 이 작을수록 assumptions violation은 치명적이다.

▶ How to solve these violations?

- ▶ (1) \Rightarrow Weighted least squares. (분산이 다르기 때문에 다른 가중치 부여)
- ▶ (2) \Rightarrow Covariance matrix (e.g., time/spatial). (y_i 's들이 independent하지 않을 때)
- ▶ (3) \Rightarrow Transformation. (Normality를 따르지 않을 때)

Nonconstant Error

- ▶ Suppose that $\epsilon_i \sim N(0, \sigma_i^2)$, $i = 1, \dots, n$ and ϵ_i 's are independent.
heteroscedasticity

\Rightarrow 등분산을 만족하지 못하면 β 의 분산이 커지게 되므로 모델의 신뢰도가 낮아짐

- ▶ Then, $\mathbf{Y} \sim MVN(\mathbf{f}, \mathbf{\Sigma})$, where $\mathbf{\Sigma} = \text{diag}(\underbrace{\sigma_1^2, \dots, \sigma_n^2}_{\text{non-constant variance}})$.

- ▶ Likelihood function:

$$L(\boldsymbol{\theta}; \mathbf{y}) = (2\pi)^{-n/2} |\mathbf{\Sigma}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{y} - \mathbf{f})^\top \mathbf{\Sigma}^{-1} (\mathbf{y} - \mathbf{f}) \right\}.$$

- ▶ MLE: For known σ_i^2 , $i = 1, \dots, n$,

$$\begin{aligned} \max_{\boldsymbol{\theta}} [l(\boldsymbol{\theta})] &\equiv \min_{\boldsymbol{\theta}} [-2l(\boldsymbol{\theta}; \mathbf{y})] \\ &\equiv \min_{\boldsymbol{\theta}} (\mathbf{y} - \mathbf{f})^\top \mathbf{\Sigma}^{-1} (\mathbf{y} - \mathbf{f}). \end{aligned}$$

Nonconstant Error

- ▶ Consider the linear regression model. That is, $\mathbf{f} = \mathbf{X}\beta$. Then MLE of β is given by

$$\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\begin{aligned} \beta_f) \quad \min_{\beta} Q(\beta) &= \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T \Sigma^{-1} (\mathbf{y} - \mathbf{X}\beta) \\ &= \min_{\beta} \mathbf{y}^T \Sigma^{-1} \mathbf{y} - 2\beta^T \mathbf{X}^T \Sigma^{-1} \mathbf{y} + \beta^T \mathbf{X}^T \Sigma^{-1} \mathbf{X} \beta \end{aligned}$$

$$\Rightarrow \left. \frac{\partial}{\partial \beta} Q(\beta) \right|_{\beta=\hat{\beta}} = -2\mathbf{X}^T \Sigma^{-1} \mathbf{y} + 2\mathbf{X}^T \Sigma^{-1} \mathbf{X} \hat{\beta} = 0$$

$$\hat{\beta} = \frac{(\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \mathbf{y}}{\text{Weighted Least Square Estimation}}$$

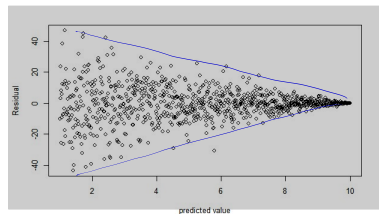
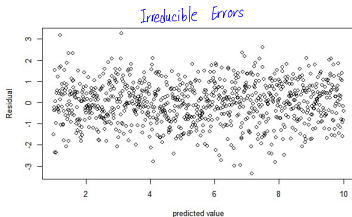
* Σ^{-1} can be called a "precision matrix"

* Σ^{-1} 의 element들이 weight가 된다 (σ_1^2 가 클수록 weight가 작아지고, σ_1^2 가 작을수록 weight가 커진다.)

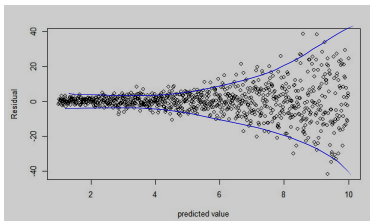
- ▶ MLE of β is the weighted least square estimator (WLSE).

Nonconstant Error with Pattern

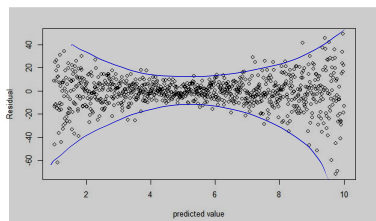
- If a residual plot shows some pattern, the variance function can be considered.



⇒ suggested transformation: coefficients with negative sign



Suggested Transformation: exponential transformation



Suggested transformation: Quadratic Transformation

Variance Function

- ▶ Variance function: $\text{Var}(\epsilon_i) = \sigma^2 g^2(\mathbf{z}_i; \boldsymbol{\theta}, \gamma)$.
 - ▶ \mathbf{z}_i : Known vector, possibly \mathbf{x}_i . \mathbf{z}_i can be $\mathbf{g}_i, \mathbf{x}_i, \mathbf{z}$
 - ▶ σ : Unknown scale parameter.
 - ▶ $g(\cdot)$: Function to be estimated by parametric or nonparametric method.
 - ▶ $\boldsymbol{\theta}$: Parameter vector of the model f .
 - ▶ γ : Parameter vector of the variance function.

Variance Function 이 mean 에 의존하는 경우도 있는데, 이러한 경우는 결국 Variance가 model f에 의존한다는 것이다.
θ가 만약 mean에 의존하지 않는다면 variance function에서 빠져나올 수 있다.

- ▶ $Y_i \sim^{indep.} N(f(\mathbf{x}_i; \boldsymbol{\theta}), \sigma^2 g^2(\mathbf{z}_i; \boldsymbol{\theta}, \gamma)), i = 1, \dots, n.$

$$y_i = f(\mathbf{x}_i; \boldsymbol{\theta}) + \sigma g(\cdot) \eta_i, \quad \eta_i : \text{pure error term} \sim N(0, 1)$$

- ▶ Examples of variance function:

- ▶ Linear pattern: $\sigma g(\mathbf{z}_i; \boldsymbol{\theta}, \gamma) = \mathbf{z}_i^\top \boldsymbol{\gamma}$.
- ▶ Exponential pattern: $\sigma^2 g^2(\mathbf{z}_i; \boldsymbol{\theta}, \gamma) = \exp(\mathbf{z}_i^\top \boldsymbol{\gamma})$.

- ▶ $\text{Var}(Y_i)$ often depends on its mean $E(Y_i)$. In that case, \mathbf{z}_i can be replaced with $\hat{y}_i = f(\mathbf{x}_i; \hat{\boldsymbol{\theta}})$.

Variance Function Estimation

Variance parameters are usually nuisance parameters

- ▶ Log likelihood function:

$$\begin{aligned} \max_{\theta, \gamma, \sigma} l(\theta, \gamma, \sigma; \mathbf{y}, \mathbf{z}) &= \max_{\theta, \gamma, \sigma} - \sum_{i=1}^n \log \{ \sigma g(\mathbf{z}_i; \theta, \gamma) \} \\ &\quad - \frac{1}{2} \sum_{i=1}^n \left\{ \frac{(y_i - f(\mathbf{x}_i; \theta))^2}{\sigma^2 g^2(\mathbf{z}_i; \theta, \gamma)} \right\}. \end{aligned}$$

- ▶ In this maximization problem, it is not easy to find θ, γ, σ simultaneously.
- ▶ Pseudolikelihood estimation:
 - ▶ To find γ and σ , it maximizes $l(\hat{\theta}, \gamma, \sigma; \mathbf{y}, \mathbf{z})$, where $\hat{\theta}$ is the MLE from $l(\theta, \hat{\gamma}, \hat{\sigma}; \mathbf{y}, \mathbf{z})$.
 - ▶ Estimations of θ and (γ, σ) are iterated until $\hat{\theta}$ is converged.

Variance Function Estimation

증명은 복잡해서 Pass !

- ▶ Residual: $r_i = y_i - f(\mathbf{x}_i; \hat{\boldsymbol{\theta}})$.
- ▶ $E(r_i^2) \approx \sigma^2 g^2(\mathbf{z}_i; \boldsymbol{\theta}, \gamma)$.
- ▶ If ϵ_i 's have normal distribution, $\text{Var}(r_i^2) \approx \sigma^4 g^4(\mathbf{z}_i; \boldsymbol{\theta}, \gamma)$.
proof is beyond the scope of this course
- ▶ Weighted estimator: γ and σ minimizing

$$\sum_{i=1}^n \frac{[r_i^2 - \sigma^2 g(\mathbf{z}_i; \gamma, \boldsymbol{\theta})]^2}{\sigma^4 g^4(\mathbf{z}_i; \gamma, \boldsymbol{\theta})}.$$

Generalized Least Squares

Algorithm

When residuals have patterns

1. Set the initial parameter vectors $\hat{\theta}$, $\hat{\gamma}$, $\hat{\sigma}$.
2. For given $\hat{\theta}$, compute squared residuals $r_i^2 = [y_i - f(x_i; \hat{\theta})]^2$.
3. Estimate the variance function parameters γ and σ by minimizing
 by profile method

$$\min_{\gamma, \sigma} \sum_{i=1}^n \frac{[r_i^2 - \sigma^2 g(\mathbf{z}_i; \gamma, \hat{\theta})]^2}{\hat{\sigma}^4 g^4(\mathbf{z}_i; \hat{\gamma}, \hat{\theta})}.$$

4. Estimate θ maximizing $l(\theta, \hat{\gamma}, \hat{\sigma}; \mathbf{y}, \mathbf{z})$.
5. Iterate Steps 2–4 until θ is converged.