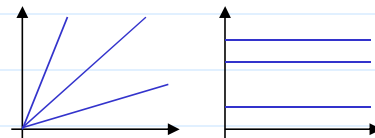


Generalized Least Square (GLS)

- model: $Y=X\beta+\varepsilon$, $E(\varepsilon)=\mathbf{0}$ and $\text{var}(\varepsilon)=\sigma^2 I \Rightarrow \varepsilon$: uncorrelated and constant variance

Q: what if $\text{var}(\varepsilon) \neq \sigma^2 I$? ε may have non-constant variance and/or are correlated, e.g.,

- time series correlation [e.g., $\varepsilon_t \sim \text{ARMA}(r,m)$]
- growth curve model, repeated measurement model [e.g., several observations taken from same person, or same unit]
- spatial correlation [e.g., data taken over contiguous geographical areas: census tracts, countries, or states in a country. Nearby areas are often much alike]
- nested errors [e.g., M sets of observations, each set from common production run, from same/common equipment, or from same survey-taker]



- Consider the case $\text{var}(\varepsilon)=\sigma^2 \Sigma$, where $\Sigma (\neq I)$ is known but σ^2 is unknown, i.e., we know the correlation and relative variance between the errors but we don't know the absolute scale

NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

- Because $\Sigma_{n \times n}$ is symmetric and positive definite, we can write $\Sigma = \mathbf{S}\mathbf{S}^T$, where \mathbf{S} is an $n \times n$ nonsingular matrix (by Cholesky or spectral decompositions)

$$Y=X\beta+\varepsilon \Rightarrow \underline{S^{-1}}Y=\underline{S^{-1}}X\beta+\underline{S^{-1}}\varepsilon \Rightarrow \underline{Y'}=\underline{X'}\beta+\underline{\varepsilon'}$$
, where

$$\underline{Y'}=\underline{S^{-1}}Y, \underline{X'}=\underline{S^{-1}}X, \underline{\varepsilon'}=\underline{S^{-1}}\varepsilon, \text{ and}$$

$$\underline{E(\varepsilon')}=\mathbf{0} \text{ and } \underline{\text{var}(\varepsilon')}=\text{var}(\underline{S^{-1}}\varepsilon)=\underline{S^{-1}}\text{var}(\varepsilon)\underline{S^{-T}}=\underline{S^{-1}}\sigma^2\underline{S}\underline{S}^T\underline{S^{-T}}=\underline{\sigma^2 I}$$

\Rightarrow For $\underline{Y'}$ and $\underline{X'}$, the assumption in ordinary least square is satisfied

- GLS: find $\underline{\beta}$ that minimize

$$\underline{\varepsilon}'^T \underline{\varepsilon}' = (\underline{Y'} - \underline{X'}\beta)^T (\underline{Y'} - \underline{X'}\beta) = (\underline{Y} - \underline{X}\beta)^T \underline{S}^{-T} \underline{S}^{-1} (\underline{Y} - \underline{X}\beta) = (\underline{Y} - \underline{X}\beta)^T \underline{\Sigma}^{-1} (\underline{Y} - \underline{X}\beta)$$

$$\Rightarrow \underline{\hat{\beta}} = (\underline{X}'^T \underline{X}')^{-1} \underline{X}'^T \underline{Y}' = (\underline{X}^T \underline{\Sigma}^{-1} \underline{X})^{-1} \underline{X}^T \underline{\Sigma}^{-1} \underline{Y}$$

$$\Rightarrow \text{var}(\underline{\hat{\beta}}) = \sigma^2 (\underline{X}'^T \underline{X}')^{-1} = \sigma^2 (\underline{X}^T \underline{\Sigma}^{-1} \underline{X})^{-1}$$

GLS is like OLS regressing

$$\underline{Y'}=\underline{S^{-1}}Y \text{ on } \underline{X'}=\underline{S^{-1}}X$$

Note1: $\underline{\varepsilon}'^T \underline{\varepsilon}'$, $\underline{\hat{\beta}}$, and $\text{var}(\underline{\hat{\beta}})$ are invariant to the choice of \underline{S} .

Note2: $\underline{\hat{\beta}} = \hat{\beta}_{\text{OLS}}$ if $\Omega[\underline{\Sigma}^{-1}X] = \Omega[X]$

- **Q:** why should not use ordinary least square when $\text{var}(\boldsymbol{\varepsilon}) \neq \sigma^2 \mathbf{I}$?

model: $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ and $\text{var}(\boldsymbol{\varepsilon}) = \sigma^2 \boldsymbol{\Sigma}$, OLS estimator $\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

$E(\hat{\boldsymbol{\beta}}_{\text{OLS}}) = \boldsymbol{\beta}$, $\text{var}(\hat{\boldsymbol{\beta}}_{\text{OLS}}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \Rightarrow$ unbiased but variance not minimized (Note: $c^T \hat{\boldsymbol{\beta}}$ is BLUE of $c^T \boldsymbol{\beta} \Rightarrow \text{var}(c^T \hat{\boldsymbol{\beta}}) \leq \text{var}(c^T \hat{\boldsymbol{\beta}}_{\text{OLS}})$)

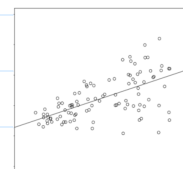
- diagnostics (residual analysis) should be applied on $\mathbf{Y}' - \mathbf{X}'\hat{\boldsymbol{\beta}} = \mathbf{S}^{-1}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{S}^{-1}\hat{\boldsymbol{\varepsilon}}$ because $\boldsymbol{\varepsilon}'$ are i.i.d. but not $\boldsymbol{\varepsilon}$
- The practical problem is that $\boldsymbol{\Sigma}$ may not be known. It's usually necessary to make some assumptions and examine the residuals to estimate $\boldsymbol{\Sigma}$ (check lab for an example, IRWLS)

❖ **Reading:** Faraway (2005, 1st ed.), 6.1 ❖ **Futher reading:** D&S, 9.2

Weighted Least Square (WLS)

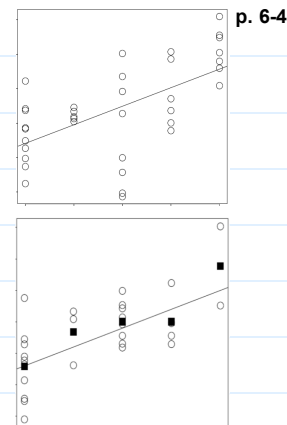
- Sometimes, the errors are uncorrelated, but have unequal variance where the form of the inequality is known ($\Rightarrow \boldsymbol{\Sigma}$ is diagonal, it's a special case of GLS), example:

- error variance proportional to a function of predictors
[e.g., $\text{var}(\varepsilon_i) = x_i^2 \sigma^2 \Rightarrow w_i = 1/x_i^2$]



NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

- data with replicates, which show a pattern of unequal variance [e.g., $\text{var}(\varepsilon_i) \approx$ sample variance of observations with same x_i ($\Rightarrow w_i = 1/(\text{sample variance})$)]
- the observed y_i 's are actually averages of several observations. [e.g., suppose y_i is the average of n_i observations, $\text{var}(\varepsilon_i) = \sigma^2/n_i$ ($\Rightarrow w_i = n_i$)]



- $\boldsymbol{\varepsilon}$: uncorrelated, but not constant variance $\Rightarrow \boldsymbol{\Sigma}$ is diagonal. Write

$$\boldsymbol{\Sigma} = \begin{pmatrix} 1/w_1 & 0 & \cdots & 0 \\ 0 & 1/w_2 & \cdots & 0 \\ \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdots & 1/w_n \end{pmatrix} \Rightarrow \boldsymbol{\Sigma}^{-1} = \begin{pmatrix} w_1 & 0 & \cdots & 0 \\ 0 & w_2 & \cdots & 0 \\ \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & \cdots & w_n \end{pmatrix}$$

where w_i 's ($\propto 1/\text{var}(\varepsilon_i)$) are called weights.

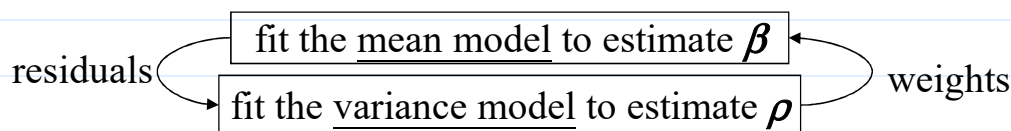
low weight \Leftrightarrow high variance; high weight \Leftrightarrow low variance

- $\mathbf{S} = \text{diag}(1/\sqrt{w_1}, \dots, 1/\sqrt{w_n})$, then $\boldsymbol{\Sigma} = \mathbf{S}\mathbf{S}^T$

\Rightarrow OLS regress $\mathbf{S}^{-1}\mathbf{Y}$ (i.e., $\sqrt{w_i} y_i$) on $\mathbf{S}^{-1}\mathbf{X}$ ($\sqrt{w_i} x_i$) (Note. the column of ones, i.e., intercept needs to be replaced with $\sqrt{w_i}$)

\Rightarrow convenient for regression package without a weighted options

- **Q:** Why observations with smaller variance should be multiplied by heavier weight? intuitive interpretation?
- iteratively re-weighted least squares (IRWLS): In all the previous examples, weights (or Σ in GLS) are assumed known. **Q:** what if $\text{var}(\varepsilon_i)$ is not completely known, what weights should we use? **Q:** where can you find the information of weights?
 - model the mean response for Y , $E(Y)=X\beta$
 - model the variance in Y , $\text{var}(Y)=f(X, \rho)$, where ρ are parameters for the variance model



- Example: $\text{var}(\varepsilon_i) = \rho_0 + \rho_1 x_{i1}$
 1. start with $w_i = 1$
 2. use weighted least square to estimate β
 3. use the residuals to estimate ρ_0 and ρ_1 , perhaps by regressing residuals² on x_{i1}
 4. re-compute the weights and go to 2. Continue until convergence

Problems: converge? how is the inference about β affected? d.f.=? ...etc

- alternative approach: jointly estimate the mean and variance parameters using likelihood based method (in R, use `gls()` function in the `nlme` library)

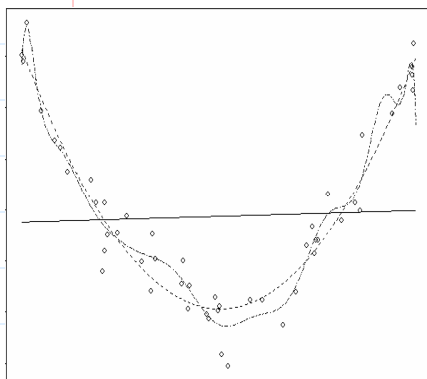
❖ **Reading:** Faraway (2005, 1st ed.), 6.2

❖ **Further reading:** D&S, 9.2

NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

Testing for Lack of Fit

- model: $Y = X\beta + \varepsilon$, **Q:** many choices of X , how can we tell the chosen X fits the data?



Q: what is “fit”? what “fit” is appropriate? example:

- data generated from the “true model” --- a 2nd-order polynomial of x
- solid line (fitted model=1st-order): not capture the quadratic pattern in data
- dashed line (fitted model=2nd-order): OK
- dotted line (fitted model=8th-order): fluctuated, fitted values and data are too close

➤ X too simple

⇒ not enough to explain the mean structure in data

⇒ lack of fit (under-fitting)

⇒ $\hat{\sigma}^2$ over-estimated

➤ X too complex

⇒ will explain the variation caused by errors, in addition to the mean structure

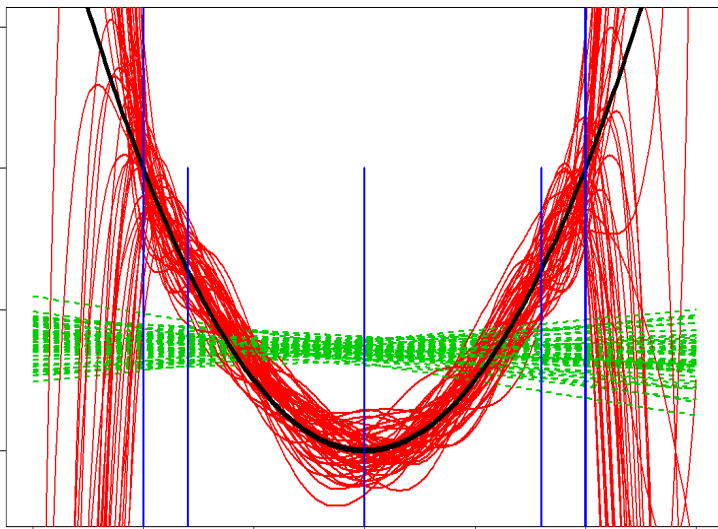
⇒ overfit the data

⇒ $\hat{\sigma}^2$ under-estimated

- **Q:** what statistic carry the information about lack of fit or overfit? Ans: $\hat{\sigma}^2$

➤ Repeat the simulation 50 times

- Black line: true model, 2nd-order polynomial
- Green lines: fitted lines under under-fitting model, 1st-order polynomials
- Red lines: fitted lines under over-fitting model, 8th-order polynomials



- Comparison
 - Under-fitting lines: small variance, large bias
 - Over-fitting lines: large variance, small bias
 - Interpolation vs. extrapolation

NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

- Note: high R^2 (a measure of “goodness of fit”) doesn't imply the model is a good fit (check lab for examples of high R^2 but lack of fit)
- graphical checking (based on residuals ← informal) (check a lab example)
- a possible testing procedure for lack of fit (\Rightarrow compare $\hat{\sigma}^2$ to σ^2)
 - when σ^2 is known
 - examples:
 - σ^2 is known from past experience
 - source of variation in ϵ is only measurement error, measuring device well understood, knowledge of measurement error inherent in an instrument or by definition
 - each y_i is an average of a very large number of observations, σ_i^2 effectively estimated (see an example in lab)
 - test for lack of fit:
 - $H_0: E(Y)=X\beta$ is correct against $H_1: E(Y)=X\beta$ is too simple
 - $(n-p)\hat{\sigma}^2 \sim \sigma^2 \chi_{n-p}^2$ if model is correct (under H_0), where n =# of observations, p =# of parameters in β
 - test statistic: $(n-p)\hat{\sigma}^2/\sigma^2 = RSS/\sigma^2$, compared with χ_{n-p}^2
 - reject if $(n-p)\hat{\sigma}^2/\sigma^2 > \chi_{n-p}^2(1-\alpha)$
 - if a lack of fit is found, a new (more complex) model is needed

➤ when σ^2 is unknown (\Rightarrow can we use data to estimate it?)

- need “model-free method” to estimate σ^2
(i.e., free of the $E(Y)=X\beta$ assumption)
- because we want to use estimated σ^2 to justify whether $E(Y)=X\beta$ is suitable,
the estimated σ^2 should have no relationship with the choice of X
- denote the estimated σ^2 under model-free method by $\hat{\sigma}_{p.e.}^2$, where p.e. stands for
“pure error”)
- usually, only possible for data with replication (**Q**: why?)

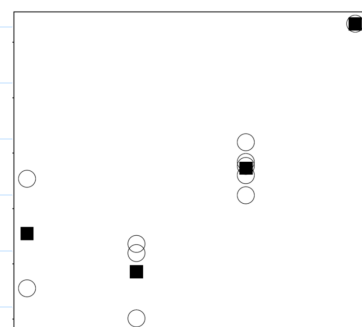
- how to estimate σ^2 (model-free)?
- $SS_{p.e.} = \sum_{\text{distinct } x} \sum_{\text{within an } x} (y_{x,i} - \bar{y}_x)^2$
- d.f. of p.e. = $\sum_{\text{distinct } x} (\# \text{ of } \text{replications} - 1)$

e.g. $(2-1)+(3-1)+(5-1)+(1-1)=7$

- $\hat{\sigma}_{p.e.}^2 = SS_{p.e.} / df_{p.e.}$

- test statistic:

(n=# of observations, p=# of parameters in β ,
RSS calculated from the model $E(Y)=X\beta$)



NTHU STAT 5410, 2022, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

| | <i>d.f.</i> | <i>SS</i> | <i>MS</i> | <i>F</i> |
|-------------|-----------------|-----------------|-----------------------------------|--|
| Residual | $n-p$ | <i>RSS</i> | | |
| Lack of fit | $n-p-df_{p.e.}$ | $RSS-SS_{p.e.}$ | $(RSS-SS_{p.e.})/(n-p-df_{p.e.})$ | ratio of <u>MS</u> (compared to $F_{n-p-df_{p.e.}, df_{p.e.}}$) |
| Pure error | $df_{p.e.}$ | $SS_{p.e.}$ | $SS_{p.e.}/df_{p.e.}$ | |

- Note: $\hat{\sigma}_{p.e.}^2$ is the estimate of σ^2 when we fit a saturated model to the data

- alternative view:

this is a comparison between the model of interest (i.e., $\omega: X\beta$)
and a saturated model (Ω , whose R^2 reaches the maximum)
that assigns a parameter to each unique combination of the
predictors \Rightarrow standard F-testing for $H_0: \omega$ v.s. $H_1: \Omega \setminus \omega$

- need replication to make the test, but it's rare in obs'nal data

- possible solution: grouping (could be questionable
 \Rightarrow different grouping schemes may cause different
conclusions)



- **Q:** what's a conservative conclusion when H_0 is accepted?
 ⇒ may not conclude $X\beta$ is the true model. We may say the true $E(Y) \approx X\beta$ on the observed data points
- **Q:** can the procedure be modified to test overfitting?
- Note that fitting is not everything
 - it often possible to fit data perfectly by introducing more effects/predictors
 - for data without replication, you can fit a model with $R^2=1$ and zero $\hat{\sigma}^2$
 - a very complex model can fit data perfectly (even exactly), but ...
 - may have no explanation (may learn nothing beyond the data itself)
 - prediction unstable
 (e.g, on region without data points, $MSE=Var+Bias^2$)
- **Q:** what is the source of variation in your data? ($X\beta$ and ϵ)
 what σ^2 is estimated (i.e., what is the source of variation in ϵ)? example:
 - replication generated from different units v.s. repeated measures of same unit
 - repeatability v.s. reproducibility in measurement system analysis

❖ **Reading:** Faraway (2005, 1st ed.), 6.3 ❖ **Futher reading:** D&S, 2.1