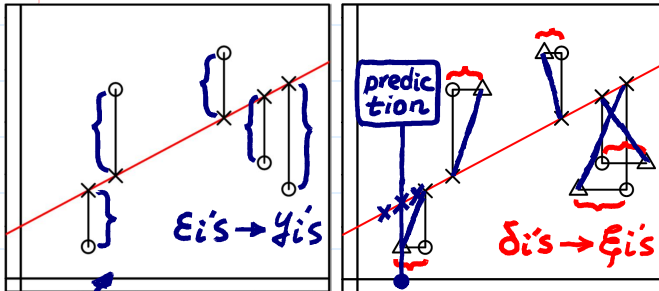# Errors in the predictor

- Recall: $Y=X\beta+\varepsilon$, where $\varepsilon$ is error that represent measurement error or unexplained variation in $Y$, and "it's assumed that $X$ are fixed values measured without error." ─random─

- **Q**: what if $X$ are measured or observed with error? (examples) ─random─ (1) sampling Data (3)DOE data (2)observational data

    Note: compare the difference between the 2 statements: "$X$ measured with errors" and "$X$ are random variables, such as in sampling model" ⇒ Both X are random variables

- **Q**: what happens if you ignore errors in $X$ and still use OLS estimator? Let us consider a simple example: ─LN p. 5-6 ~ 7. $E(y_\xi)=\beta_0+\beta_1\xi$    but, different causes of their randomness

    

    $X$: "fine" value, no error in predictor and response, but not observed
    ○: observations with error only in response, but not error in predictor
    △: observations with error in both response and predictor

    ε_i's → y_i's        δ_i's → ξ_i's

- **Q**: Is minimizing *RSS* (i.e., OLS) still reasonable for △ data? **No.**

- a statistical model for △ data: $E(y_\xi)$        not observed        not a linear model
    ideal "true" relationship is $\eta_i=\beta_0+\beta_1\xi_i$,        X: $(\xi_i, \eta_i)$
    but observe        $y_i=\eta_i+\varepsilon_i$ and $x_i=\xi_i+\delta_i$,        ○: $(\xi_i, y_i)$
    where $\varepsilon$ and $\delta$ are errors of response and predictor respectively, i.e.,        △: $(x_i, y_i)$
    [assume independent]  [observed] $y_i=\beta_0+\beta_1\xi_i+\varepsilon_i=\beta_0+\beta_1 x_i+(\varepsilon_i-\beta_1\delta_i)$    observed  $\varepsilon_i^*$

---

- **Q**: what problem if we use ordinary least square to estimate $\beta_1$ in the model?        c.f. $Var(\varepsilon_i^*) = \sigma^2 + \beta_1^2\sigma_\delta^2$

    - Let's assume $E(\varepsilon_i)=E(\delta_i)=0$, $var(\varepsilon_i)=\sigma^2$, $var(\delta_i)=\sigma_\delta^2$ and $cov(\varepsilon,\delta)=0$

    - ξ_i's are random in sampling model
    Let $\sigma_\xi^2=\Sigma(\xi_i-\bar\xi)^2/n$ (Note: when $\xi_i$'s are not random, we could regard it as a measure of the spread of the predictor), $\sigma_{\xi\delta}=cov(\xi,\delta)$ and assume $cov(\xi,\varepsilon)=0$

    - the OLS estimator of $\beta_1$ is:        cov(x,y)/Var(x)
    $\hat\rho\frac{\hat\sigma_y}{\hat\sigma_x}=\hat\beta_1=[\Sigma(x_i-\bar x)(y_i-\bar y)]/\Sigma(x_i-\bar x)^2$    $\frac{cov(\xi+\delta, \beta_1\xi+\varepsilon)}{Var(\xi+\delta)}$
    LN p. 3-7

    - after some calculation, we can write
    (exercise) → $E(\hat\beta_1)\approx\beta_1\times[(\sigma_\xi^2+\sigma_{\xi\delta})/(\sigma_\xi^2+\sigma_\delta^2+2\sigma_{\xi\delta})]$    unbiased if $\sigma_{\xi\delta}=-\sigma_\delta^2$

    - if no relation between $\xi$ and $\delta$ (i.e., $cor(\xi,\delta)=0$),
    $E(\hat\beta_1)\approx\beta_1\times[\sigma_\xi^2/(\sigma_\xi^2+\sigma_\delta^2)]=\beta_1\times[1/(1+\sigma_\delta^2/\sigma_\xi^2)]$ ⇒ $\hat\beta_1$ is biased

    - typically, bias in $\hat\beta_1$ is towards zero    ≤ 1        if X & δ are uncorrelated, $\sigma_{\xi\delta}=cov(\xi,\delta)=cov(x-\delta,\delta)=-\sigma_\delta^2$

    - size of the bias depends mainly on the ratio $\sigma_\delta^2/\sigma_\xi^2$ (i.e., variability in the errors of predictor relative to the spread of predictor) (**Q**: why reasonable?)

    How? ①reduce $\sigma_\delta^2$ ②increase $\sigma_\xi^2$
    - ratio is small ⇒ no worry
    - ratio is large, $|\hat\beta_1|$ is underestimated ⇒ use *measurement error model*

- For multiple predictors, the usual effect of measurement errors on predictors is to bias the estimator of $\beta$ in the direction of zero    But, $Var(\hat Y)$↑ ∵ error in X

- Prediction is not biased since future $X$ will also be measured with errors. So, model for prediction should be built on $X$'s measurement with error.    Note. $y=\beta_0+\beta_1 x+\varepsilon^*$

❖ **Reading**: Faraway (1st ed.), 5.1; W, 4.6.3    ❖ **Further reading**: D&S, 3.4, 9.7

# Collinearity ← Recall. Identifiability (LNp.5-11~12)

p. 9-9

- collinearity: predictors are (linearly) related to each other    model: $y = \sum_j \beta_j g_j(\underline{x}) + \varepsilon$

  ➢ $X^TX$ is singular $\Rightarrow$ some predictors are linear combinations of others    unidentifiable
  $\Rightarrow$ (exact) collinearity $\Rightarrow$ no unique estimate of $\beta$  $\exists a_1, \cdots, a_p$ s.t. $\sum_{j=1}^{p} a_j g_j(\underline{x}_i) = \underline{0}$, $\forall i$

  ➢ $X^TX$ close to singular $\Rightarrow$ close to linear dependent among some predictors
  $\Rightarrow$ (approximate) collinearity or multicollinearity  $\exists a_1, \cdots, a_p$ s.t. $\sum_{j=1}^{p} a_j g_j(\underline{x}_i) \approx 0$

- effect of collinearity: $\overbrace{\text{(exercise by Schur complement)}}$ $\leftarrow (X^TX)^{-1}_{jj} = [\underline{g_j^T g_j} - \underline{g_j^T X_{(-j)}(X_{(-j)}^T X_{(-j)})^{-1} X_{(-j)}^T g_j}]^{-1} \leftarrow \underline{RSS_j^{-1}}$

  ➢ estimated effects are unstable (can change magnitude or sign depending on the
  $\boxed{S_j \leftarrow TSS_j}$ other predictors in the model) $\Rightarrow$ interpretation of estimated coefficients difficult

  ➢ cause numerical problem in estimating $\beta$ and associated quantities $\leftarrow$ calculate $(X^TX)^{-1}$

  $\boxed{\begin{array}{c}(X^TX)^{-1}_{jj}\\ \times \sigma^2\end{array}}$ ➢ var($\hat{\beta}_j$) = $\sigma^2$ ($1/(1-R_j^2)$) ($1/S_j$), where $S_j = \sum_i (g_{ij} - \bar{g}_j)^2$ and $R_j^2$ is the coefficient of determination obtained from regressing $g_j$ on all other predictors $\Rightarrow$ when $R_j^2 \approx 1$, var($\hat{\beta}_j$) large $\Rightarrow$ $t$-test may fail to reveal significance, i.e., miss important $g_j$

  ➢ *variance inflation factor*: $VIF_j = 1/(1-R_j^2) \Rightarrow$ when $S_j$ is fixed, $VIF_j$ represents the increase in variance due to the collinearity (e.g., interpret $VIF_j = 16$?)

- detection of collinearity: ← compared to the case of orthogonality (i.e., $R_j^2 = 0$)

  ➢ examine correlations between predictors, i.e., cor($g_k, g_j$) ← from $X^TX$
  $\Rightarrow$ any values close to 1 or –1 reveal *pairwise* correlation

  $\boxed{\text{s.e.}(\hat{\beta}_j) \approx \sqrt{16} = 4 \text{ times larger than being orthogonal.}}$

  ➢ for each $g_j$, regress $g_j$ on all other predictors and compute $R_j^2$ or $VIF_j$
  $\Rightarrow$ $R_j^2$ close to one or $VIF_j$ much larger than one indicate a problem of collinearity

---

p. 9-10

  ➢ examine eigenvalues, $\lambda_1 \geq \ldots \geq \lambda_p$, of $X^TX \Rightarrow$ small eigenvalues indicate a problem
  $\boxed{LNp.5.12}$
  ↳ eigenvector $(a_1, \cdots, a_p)$

  - condition number: $k = (\lambda_1/\lambda_p)^{1/2}$       Then, $a_1 g_1 + \cdots + a_p g_p \approx \underline{0}$
  - rough rule: $k > 30$ is considered large $\Rightarrow \lambda_1/\lambda_p > 900$
  - for each $i$, $(\lambda_1/\lambda_i)^{1/2}$ are worth considering $\Rightarrow$ there may exist more than one linear combination relationship between predictors
  - eigenvectors of small eigenvalues indicate possible source of collinearity

- how to deal with collinearity:

  ➢ identify the cause of collinearity in data ← $\boxed{\text{Check LNp.5-11~12}}$ ↱ explain why collinearity occurs, not only to detect whether it occurs.

  check 1st Note in LNp9-6 ➢ amputate some predictors if you can --- remember that collinearity happens because too many variables try to do the *same job* of explaining the response ← cf.

  ➢ do not conclude the predictors we drop have nothing to do with the response ← cf.

  ➢ techniques such as principle component regression, ridge regression, partial least squares, …, may help ↱ to reduce the impact of collinearity, e.g., PC use linear combinations of $g_j$'s

❖ **Reading**: Faraway (1st ed.), 5.3; W, 10.1
❖ **Further reading**: D&S, 16.1, 16.4, 16.5

every column of $Z$ is a linear combination of the columns of $X$

主成分 ⟶ **Principal components**

- Recall: $Y = X\beta + \varepsilon$. If $X$ is orthogonal (i.e., $X^TX$ is a diagonal matrix), then estimation, testing, and parameter interpretation are greatly simplified.

  $\boxed{\beta, \beta' \text{ are different parameters}}$

- idea: For non-orthogonal $X$, replace $Y = X\beta + \varepsilon$ by $Y = Z\beta' + \varepsilon$, where $Z$ is a linear combinations of $X$ (i.e., $Z_{n \times q} = X_{n \times p} U_{p \times q}$, $p \geq q$) and $Z$ is orthogonal ($Z^TZ$ is diagonal)

- orthogonality
  Recall: (LNp. 4-5)
  bivariate
  normal    $\bar{g}=0$
  orthogonal $g_j, g_k$
  $\Rightarrow cor(g_j, g_k)=0$

model matrix $=[\mathbb{1} \ g_1 \ g_2 \cdots g_p]=[\mathbb{1} \ X]$

1st coordinates



$g_1, g_2$ independent (cor=0) — $g_1, g_2$ correlated (cor $\neq$ 0)
2nd coordinates

- $\underline{Z}_{n\times q} = \underline{X}_{n\times p}\,\underline{U}_{p\times q}$, $p\geq q$, e.g., take a look of the first column of $Z$

rotation matrix — 1st row of $X$ "●" — projection
2nd column of $U$ | 1st column of $U$

$$Z_1 \equiv \begin{bmatrix} z_{11} \\ z_{21} \\ \cdots \\ z_{n1} \end{bmatrix} \xleftarrow{\text{inner product}} \begin{bmatrix} g_{11} & g_{12} & g_{1p} \\ g_{21} & g_{22} & g_{2p} \\ \cdots & \cdots & \cdots \\ g_{n1} & g_{n2} & g_{np} \end{bmatrix}$$

$u_{11}$ … $u_{21}$ … $u_{p1}$ — 1st column of $U$ (unit vector)

$(x-\bar{x})'S^{-1}(x-\bar{x}) = c^2$
2nd coordinate
$z_{12}$
$z_{11}$ 1st coordinate
row of $X$
new (0,0)
original (0,0)

$(g_{11}, g_{12}, \ldots, g_{1p}) \xrightarrow{\text{rotation}} (z_{11}, z_{12}, \ldots, z_{1q})$ 1st row of $Z$ "●"

- see graph, in which $z_1$ is the projection of points on the
  direction $u_1$; $z_2$ is the projection of points on the direction $u_2$
  cf. → model selection.

cf. graph in LNp. 3-3

PC

- concept of *dimension reduction*: → [1. smaller dimension, better; 2. important information should be kept]

  - **Q**: take a look of the graph, the points are of 1-dim or of 2-dim? model selection
    $\Rightarrow$ very similar to a line $\Rightarrow$ high correlation $\Rightarrow$ data is 2-dim, but close to 1-dim

  - replace large number of columns in $X$ with small number of columns in $Z$ (usually $\leq p$)
    $\Rightarrow$ simpler model, especially useful (1) when few linear combinations of $X$ are
    enough to represent the variation in $X$; (2) when $p > n$ ← unidentifiable

---

- principal component (PC): $Z_j = X\,U_j \Leftarrow Z = X\,U$
  $z_j^T z_j = U_j^T(X^TX)U_j = \lambda_j \|U_j\|^2 = \lambda_j$

  - transform $X$ to $Z$ which is orthogonal, but how?
    $Z_j, Z_k$ uncorrelated, if $E(Z_j)=E(Z_k)=0$

  - find $U$ such that ⭐ $Z^TZ$ is diagonal, i.e., $Z^TZ = \text{diag}(\lambda_1,...,\lambda_p)$, where $\lambda_1 \geq ... \geq \lambda_p \geq 0$

  - since $Z^TZ = U^T(X^TX)U$, to make $Z^TZ$ diagonal, we can choose columns of $U$ are
    (symmetric) orthogonal eigenvectors of $X^TX$, then the $\lambda_1, \lambda_2, ..., \lambda_p$ are eigenvalues of $X^TX$

    - let $U_j$ and $\lambda_j$ be the $j$-th eigenvector and eigenvalue of $X^TX$, then $(X^TX)U_j = \lambda_j U_j$
    - $U_k^T U_j = 0$ for $k \neq j$ and $\|U_j\|=1$ for all $j$
    - $(Z^TZ)_{kj}$: $U_k^T(X^TX)U_j = \lambda_j U_k^T U_j$, which equals 0 if $k \neq j$ and equals $\lambda_j$ if $k=j$

  - $Z_1$ (=1st column of $Z$) is called 1st principal component (PC),
    $Z_2$ (=2nd column of $Z$) is called 2nd principal component (PC), ...
    projection

  - another way to look at it: $z_j^T z_j = (Z^TZ)_{jj} = \lambda_j$

$u_1$ $u_2$
$\lambda_2$ $\lambda_1$
(0,0)

$U_1$: unit vector
$\|Z_j\|^2$ $=z_j^T z_j$ $=(XU_j)^T$ $(XU_j)$ $=U_j^T X^T X U_j$ $\Rightarrow$ quadratic form

- $Z_1$ = linear combination of columns of $X$ that has maximum length$^2$, i.e, maximizing $\Sigma z_{i1}^2$ (variation of $Z_1$)
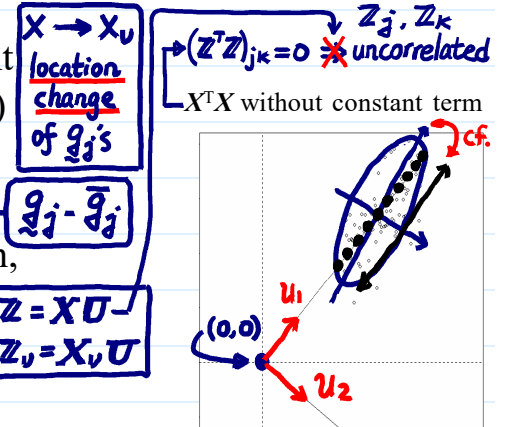- $Z_2$ = linear combination of columns of $X$ that is orthogonal to $Z_1$ and has maximum length$^2$ ⊕ $U_2 \perp U_1$ ($\bar{g}_j = 0 \ \forall j$)
- $Z_3$ = linear combination of columns $X$ that is orthogonal to $Z_1, Z_2$ and has maximum length$^2$ ⊕ [$U_3 \perp U_1$; $U_3 \perp U_2$]
- …

- some properties:
  - ⭐ $U^T U = I_{q\times q}$ ← $U$: an orthogonal matrix
  - zero eigenvalue $\Rightarrow$ unidentifiable
  - $\lambda_j$ = length$^2$ of $Z_j$ = $\Sigma_i z_{ij}^2$ [note: when $E(X_j)=0 \Rightarrow E(Z_j)=0 \Rightarrow \lambda_j \propto var(Z_j)$]
    sample mean — sample variance
  - $\lambda_1 + ... + \lambda_p = \text{tr}(X^TX) = \Sigma_j(\text{length}^2 \text{ of } X_j)$
    [note: when $E(X_j)=0, \ \forall j$ — sum of variations of different units
    $\text{tr}(Z^TZ)$ $\lambda_1 + ... + \lambda_p \propto \Sigma_j var(X_j)$: total variation of $X$
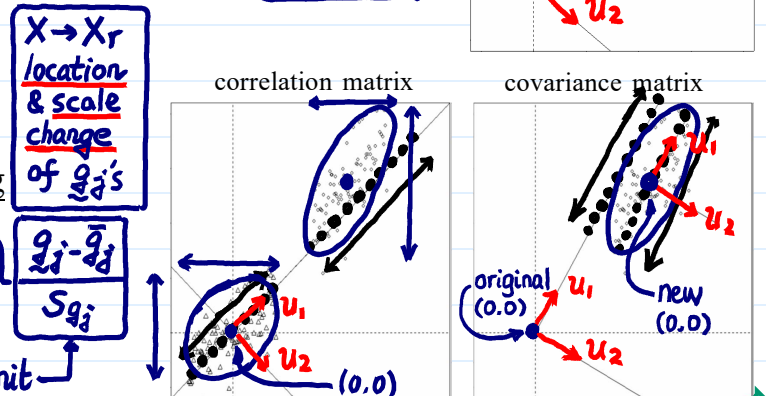  - $\lambda_j / (\lambda_1 + ... + \lambda_p)$ = proportion of total variation explained by the $j$th PC

➢ **Q**: how to interpret $Z_1, Z_2,..., Z_p$? Ans: compare the coefficients in eigenvectors
- Ex 1: $Z_1$=0.46GNP+0.32UnEm+0.46POP+0.46Year+... ⇒ hard to give meaning
  *average*
- Ex 2: $Z_1$=0.63(hw1)+0.57(hw2)+0.52(hw3) ∝ average homework scores;
  $Z_2$ = 0.67(hw1)+0.08(hw2)−0.75(hw3) ∝ difference between hw 1 and 3 scores
  *≈0*

• variation on principal component regression
  ➢ use $X^T X$ with/without constant term (without constant
    term ⇒ PC's may not be orthogonal to constant term)

    > $X \to X_v$ location change of $\underline{g}_j$'s
    > $\underline{g}_j - \bar{\underline{g}}_j$

  ➢ use covariance matrix of $\underline{X}$ (without constant term),
    i.e., $X_v^T X_v/(n-1)$ where $X_v$ is formed by centering
    each $\underline{g}_j$, to find eigenvectors $U$ and eigenvalues. Then,
    $\lambda_j$ = var($z_j$). The transformation $U$ can be applied on
    $\underline{X}$ or $X_v$    [PC's are orthogonal to constant term if
    transformation is applied on $X_v$]

    > $Z = XU$
    > $Z_v = X_v U$

    $z_j, z_k$ → $(Z^T Z)_{jk}=0$ ✗ uncorrelated

    $X^T X$ without constant term

    
    cf. $u_1$ (0,0) $u_2$

  ➢ use correlation matrix of $X$ (without
    constant term), i.e, $X_r^T X_r/(n-1)$,
    where $X_r$ is formed by standardizing
    each $\underline{g}_j$. To make sense, the
    transformation should be applied on
    $X_r$ .Then, $\lambda_j$ = var($z_j$) and PC's are
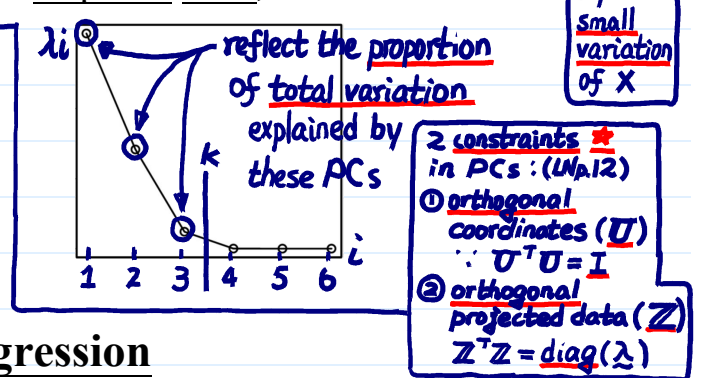    orthogonal to constant term

    > $Z_r = X_r U$

    > $X \to X_r$ location & scale change of $\underline{g}_j$'s
    > $\dfrac{\underline{g}_j - \bar{\underline{g}}_j}{S_{\underline{g}_j}}$
    > *free of unit*

    correlation matrix              covariance matrix

    
    $u_1$ $u_2$ (0,0)     original (0,0) $u_1$   new (0,0) $u_2$

---

⬅ • Notes:
  ➢ interpretation is a problem --- little is gained if
    principal components are not interpretable ⬅

    *Note. In the development of PCs, the criteria (LNp.12) do not consider interpretation*

  ➢ how many principal components are worth considering? plot $\lambda_i$, often the plot has
    a noticeable "elbow" --- the point, say $k$, at which further eigenvalues are
    negligible in size compared to the earlier ones ⇒ $(\lambda_1+...+\lambda_k)/(\lambda_1+...+\lambda_p)$ =
    proportion of total variation explained by the first $k$ principal components

  ➢ principal components do not use information from the response in
    dimension reduction. It is possible that a lesser principal component is actually
    very important in explaining/predicting the response. Dimension-reduction
    methods that utilize information about the response exist, such as
    - partial least square ⬅ *relax one* ★
    - sliced inverse regression (SIR)
    - principal Hessian directions (pHd)
    - projection pursuit regression
    - canonical correlation analysis
    - LASSO ⬅

    *Recall criteria in LNp 12*

    *explain small variation of X*

    
    $\lambda_i$   *reflect the proportion of total variation explained by these PCs*   $k$   1 2 3 4 5 6   $i$

    *2 constraints ★ in PCs: (LNp.12)*
    ① *orthogonal coordinates ($U$)* ∵ $U^T U = I$
    ② *orthogonal projected data ($Z$)* $Z^T Z = diag(\underline{\lambda})$

  ❖ **Reading**: Faraway (1st ed.), 9.1

    *cf.*  *other better estimator than OLS?* ⬅

### Ridge regression

• **Q**: what is the problem? strong collinearity (i.e., $X^T X$ close to singular) causes (1)
  numerical problem in calculating $(X^T X)^{-1}$; (2) $\hat{\beta}$ unstable; (3) large variance in $\hat{\beta}$
  *OLS estimator*

- ridge estimator: a method of <u>combating</u> <u>strong collinearity</u> [<u>Note</u>: It would be <u>better</u> to find out <u>how collinearity occurs</u> before doing <u>ridge regression</u>.]

  *can compare* $\hat{\gamma}$

  *different range of an X in different data*

  ➢ <u>centering and scaling predictors</u>: $X \to F$, i.e., $F^T F$ = <u>correlation matrix</u> of $X$

  (**Q**: <u>why</u>?), and <u>centering response</u>: $Y \to Z$, i.e., $\underline{Z} = \underline{Y - \bar{Y}}$,

  *Standardization*

  *∵ + $\lambda I$*

  *OLS est'or of $\gamma$*
  $= (F^T F)^{-1} F^T Z$
  *unstable if there exists strong collinearity*

  $$y = \beta_0 + \beta_1 x_1 + ... + \beta_p x_p + \varepsilon \implies Y = X\beta + \varepsilon$$
  $$z = \boxed{gone}\ \gamma_1 f_1 + ... + \gamma_p f_p + \varepsilon \implies Z = F\gamma + \varepsilon$$

  *Recall Location & Scale change (LNp.8-2)*

  Note: $\underline{\beta_i = \gamma_i / sd_i}$, where $\underline{sd_i}$ is the <u>sample standard deviation</u> of $\underline{x_i}$, $i=1, ..., p$

  ➢ <u>ridge estimator</u>: for $\lambda > 0$, $\hat{\gamma}_\lambda = (F^T F + \underline{\lambda I})^{-1} F^T Z = (F^T F + \lambda I)^{-1} F^T Y$   [note: $F^T \mathbf{1} = \mathbf{0}$]

  *$F^T F$ & $F^T F + \lambda I$ have same eigenvectors*

  - $\lambda = 0 \implies \underline{\hat{\gamma}_0}$ is the <u>OLS estimator</u> and $\lambda \to \infty \implies \hat{\gamma}_\infty = \mathbf{0}$ *no intercept*

    *OLS est'or* $\gamma_2$   $\lambda \uparrow$   $\gamma_1$   Q → *shrinkage*

  - for an <u>eigenvector</u> $\underline{u_i}$ of $F^T F$ and its corresponding eigenvalue $\lambda_i$, $(F^T F + \lambda I)u_i = (\lambda_i + \lambda)u_i \implies u_i$ is an <u>eigenvector</u> of $(F^T F + \lambda I)$ with corresponding eigenvalue $\underline{\lambda_i + \lambda}\ (>\lambda_i)$ ← *strong collinearity* ⟺ *some $\lambda_i$'s $\approx 0$*
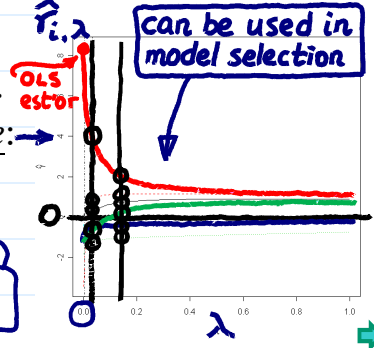
  - <u>ridge estimator</u> can <u>remedy</u> the <u>problems</u> caused by <u>strong collinearity</u>

  (Δ) ➢ <u>how</u> to <u>choose</u> an <u>appropriate</u> $\lambda$? → *criteria, e.g., crossvalidation, ...*

  *$F^T F = U\Lambda U^T$. $F^T F + \lambda I = U(\Lambda + \lambda I)U^T$ ∵ $UU^T = I$*

  There exists various <u>methods</u> *automatically* choosing a $\lambda$. However, the <u>most popular method</u> is through *ridge trace*: →

  *why?* →

  plot $\hat{\gamma}_\lambda$ against $\underline{\lambda}$

  Find a <u>minimum value</u> of $\underline{\lambda}$ (usually chosen in [0, 1]) *cf* after which $\hat{\gamma}_\lambda$ are <u>moderately stable</u>. $\sum_{i=1}^{p} \lambda_i = trace(F^T F) = p$ *average of $\lambda_i$'s = 1*

  $\hat{\gamma}_{i,\lambda}$ *can be used in model selection*

  *OLS est'or*