

Regression

Predictor: Variable $X_1 \dots X_d$ ($n \times 1$ feature vectors)

Response: RV Y ($n \times 1$)

- Conditional expectation of Y given values $x_1 \dots x_d$ of var $X_1 \dots X_d$ is called the **regression function** of Y on $X_1 \dots X_d$, or simply the regression of Y on $X_1 \dots X_d$. In symbols, $E(Y|x_1 \dots x_d)$

Assumption of LR. (General Linear Model)

1. Predictor is known (for multiple reg. requires $\text{rank}(X) = d$)
2. Normality: conditional distr of Y_i is normal.
3. Linear Mean: \exists parameter β_{d+1} , $E(Y|x_1 \dots x_d) = \sum_i \beta_i x_i$
4. Common Variance (homoscedasticity): $\exists \sigma^2$, $\text{Var}(Y_i|x_1 \dots x_d) = \sigma^2$, $i=1 \dots n$.
5. Independence: Given $x_1 \dots x_d$, RV $Y_1 \dots Y_n$ are independent.

Given the above assumptions, the conditioned joint distr of $Y_1 \dots Y_n$ given vector $x_1 \dots x_d$, parameters β , σ^2 is:

$$f(Y|X, \beta, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \sum (Y_i - \sum_j \beta_j x_j)^2\right)$$

Formulation of LR

- Simple LR

$$y = \beta_0 + \beta_1 X + \varepsilon$$

- Multiple LR

$$y = X_0 \beta_0 + \dots + X_d \beta_d + \varepsilon \quad X_i : n \times 1$$

Note, constant term is generalized with $X_0 = \mathbf{1}$.

Matrix notation: $y = X\beta + \epsilon$, $X: n \times d$, $\beta: d \times 1$

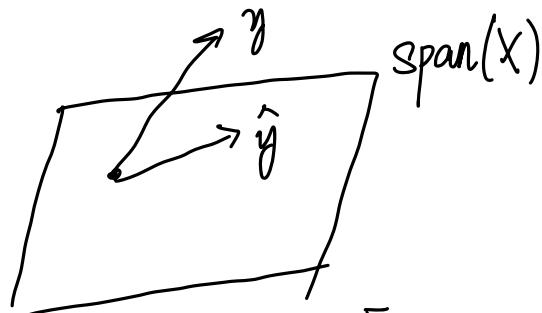
Solution of LR

Under L2 loss: $L(\beta) = \|y - X\beta\|_2^2$ $\det(X^T X) \neq 0$

$$\hat{\beta} = \begin{cases} (X^T X)^{-1} X^T y & \text{when } \det(X) \neq 0 \text{ (invertible)} \\ V \Sigma^+ U^* y & \text{when } \det(X) = 0, \text{ SVD: } X = U \Sigma V^* \\ R^T Q^T y & \text{when QR: } X = QR \end{cases}$$

Geometrical Perspective of OLS

$$y = X\beta, \quad \text{OLS } \hat{\beta} = X^T y$$



Note: $X X^+$ is a linear orthogonal operator that projects any vector onto $\text{span}(X)$.

From OLS, $\hat{y} = X\hat{\beta} = X X^+ y$, which is projection of y .

Gauss - Markov Theorem

Least square solution is of least variance among all linear unbiased estimators of parameter β .

Formal proof:

Let $\tilde{\beta} = Cy$ be an arbitrary unbiased linear estimator of β .

Choose D , s.t., $C = X^+ + D$

$$\begin{aligned}\underbrace{\mathbb{E}(\tilde{\beta})}_{\text{unbiasedness}} &= \beta : \mathbb{E}(\tilde{\beta}) = \mathbb{E}(Cy) = \mathbb{E}(X^+ + D)(X\beta + \varepsilon) \\ &= \beta + DX\beta + (X^+ + D)\mathbb{E}(\varepsilon) \\ &= (I + DX)\beta \quad \underbrace{= 0}\end{aligned}$$

Thus, in order for $(I + DX)\beta = \beta$, $DX = 0$.

$$\begin{aligned}\text{Var}(\tilde{\beta}) &= \text{Var}(Cy) = C \text{Var}(y) C^T = \sigma^2 CC^T \\ &= \sigma^2 [(X^+ + D)(X^+ + D)^T] = \sigma^2 [X^+ X^{+T} + X^+ D^T + DX^{+T} + DD^T] \\ &= \sigma^2 [(X^+ X)^{-1} X^T X (X^+ X)^{-1 T} + (X^+ X)^{-1} X^+ D^T + DX(X^+ X)^{-1 T} + DD^T] \\ &\quad (\text{DX} = 0) \\ &= \sigma^2 [(X^+ X)^{-1} + DD^T] = \text{Var}(\hat{\beta}) + \sigma^2 DD^T\end{aligned}$$

Since DD^T p.s.d, $\text{Var}(\tilde{\beta}) \geq \text{Var}(\hat{\beta})$. □

Alt. explanation for naive case:

We will use bias-variance tradeoff for estimators.

Let $\hat{\theta}$ be an estimator of θ .

Note that, we should not treat this as a Bayesian estimation where θ is RV, but rather use frequentist view that is conditioned on θ , i.e., θ is regarded as constant.

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= \mathbb{E}(\hat{\theta} - \theta)^2 = \mathbb{E}(\hat{\theta}^2) - 2\mathbb{E}(\hat{\theta})\theta + \theta^2 \\ &= \mathbb{E}(\hat{\theta}^2) - \mathbb{E}^2(\hat{\theta}) + \mathbb{E}^2(\hat{\theta}) - 2\mathbb{E}(\hat{\theta})\theta + \theta^2 \\ &= \text{Var}(\hat{\theta}) + [\mathbb{E}(\hat{\theta}) - \theta]^2 \end{aligned}$$

Variance of estimator || bias squared

We will showcase Gauss-Markov for a linear combination $\theta = \alpha^T \beta$.
(Here, we can regard both α and β as constant)

Now, least square estimator is $\hat{\theta} = \hat{\alpha}^T \hat{\beta} = \alpha^T X^+ y$

It's unbiased if we assume the linear model is correct:

$$\mathbb{E}(\hat{\alpha}^T \hat{\beta}) = \alpha^T X^+ \mathbb{E}(y) = \alpha^T X^+ X \beta = \alpha^T \beta$$

For an arbitrary linear estimator wrt y , $C^T y = \hat{\theta}$.

By bias-var decomposition, $\text{MSE}(\hat{\theta}) = \text{Var}(\hat{\theta}) + \text{bias}^2(\hat{\theta})$

Thus for unbiased estimators, $\text{Var}(\hat{\theta}) = \text{MSE}(\hat{\theta})$

which means the least square estimator which minimizes MSE also minimizes Variance.

Estimators of β and σ^2

treat RSS as constant. min f gives $\hat{\beta}^2$.
 min RSS, or least square gives $\hat{\beta}$.

MLE: comes from data likelihood function $f(Y|f, \theta) = C \left(\frac{1}{\sigma^2}\right)^{\frac{n}{2}} e^{-\frac{1}{2\sigma^2} \text{RSS}}$

$$\begin{aligned}\hat{\beta}^1 &= \text{least square estimate} \\ \hat{\sigma}^2 &= \frac{1}{n} \sum (Y_i - (\sum \beta_j X_j))^2 = \frac{1}{n} \text{RSS}\end{aligned}\quad \boxed{\hat{\beta}, \hat{\sigma}^2 \text{ are still RVs}}$$

Unbiased:

$$\beta' = \hat{\beta} \quad , \quad \sigma'^2 = \frac{1}{n-d} \text{RSS} = \frac{n}{n-d} \hat{\sigma}^2$$

Distribution of $\hat{\beta}_i$ (Simple LR)

Simple LR, $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)^T$, \vec{x} fixed, \vec{Y} RV.

$$\hat{\beta}_1 = \frac{SXY}{SXX} \quad , \quad \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$$

Intuition: slope need to be of form Y/X

From the above definition of $\hat{\beta}_i$, taking E and Var of $\hat{\beta}_i$ and using assumption $E(Y) = \beta_0 + \beta_1 x$, $\text{Var}(Y) = \sigma^2$,

We get distr of $\hat{\beta}_i$:

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{SXX}\right), \quad \hat{\beta}_0 \sim N\left(\beta_0, \sigma^2\left(\frac{1}{n} + \frac{\bar{X}^2}{SXX}\right)\right),$$

$$\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\bar{X}\sigma^2}{SXX}$$

- Estimator of $\text{Var}(\hat{\beta}_i)$ can be obtained by substituting σ with estimators of σ , e.g., $\hat{\sigma}$, σ' .

Joint Distribution of $\hat{\beta}$, $\hat{\sigma}^2$ (Simple LR)

- Joint distr of $(\hat{\beta}_0, \hat{\beta}_1)$ is bivariate normal w/ μ, Σ above.

- For $n \geq 3$, $\hat{\sigma}^2$ is independent of $(\hat{\beta}_0, \hat{\beta}_1)$

RSS is also independent of $\hat{\beta}$, because $\hat{\sigma}^2 = \frac{1}{n} \text{RSS}$

- $n \hat{\sigma}^2 / \sigma^2 = \frac{\text{RSS}}{\sigma^2}$ has χ^2 distr w/ $n-2$ dof.

$$\text{Intuition: } \frac{n \hat{\sigma}^2}{\sigma^2} = \frac{\sum (Y_i - X_i^T \hat{\beta})^2}{\sigma^2}$$

We use $\hat{\beta}$ to estimate β as it's unknown.

Suppose we have known β , and use that in the above eq,

$$\Rightarrow \frac{n \hat{\sigma}^2}{\sigma^2} = \frac{\sum (Y_i - X_i^T \beta)^2}{\sigma^2}$$

By assumption of LR, $Y_i \sim N(X_i^T \beta, \sigma^2)$.

Thus $(Y_i - X_i^T \beta)^2 / \sigma^2 \sim \chi^2$ w/ 1 dof.

$$n \hat{\sigma}^2 / \sigma^2 = \sum (Y_i - X_i^T \beta)^2 / \sigma^2 = \sum (\chi^2 \text{ w/ 1 dof}) = \chi^2 \text{ w/ } n \text{ dof.}$$

Now by estimating β with $\hat{\beta}$, we removed

$\begin{cases} 2 \text{ dof } (\hat{\beta}_0, \hat{\beta}_1) \text{ for simple LR from } \chi_n^2 \\ d \text{ dof } (\hat{\beta}_{d+1}) \text{ for multiple LR} \end{cases}$

or, $n \hat{\sigma}^2 / \sigma^2 \sim \chi^2$ w/ $n-2$ dof for simple LR.

Draft Proof:

Let's construct orthogonal matrix $A_{n \times n}$:

1st column of A : $a_{1j} = \frac{1}{\sqrt{n}}$

2nd column of A : $a_{2j} = \frac{1}{\sqrt{S_{xx}}} (x_j - \bar{x})$

Pad other columns with Gram-Schmidt.

Now, let $Z = AY$ ($Z_{n \times 1}, Y_{n \times 1}$)

It's possible to show:

- Z_1, \dots, Z_n are independent and has normal distr w/ σ^2 .

Now, since $Z_1 = \sum a_{1j} Y_j = \sqrt{n} \bar{Y}$, $\hat{\beta}_0 = \bar{Y} - \bar{x} \hat{\beta}_1$,

$$Z_1 = \sqrt{n} (\hat{\beta}_0 + \bar{x} \hat{\beta}_1).$$

Similarly, $Z_2 = \sum a_{2j} Y_j = \frac{1}{\sqrt{S_{xx}}} \sum (x_j - \bar{x}) Y_j = \sqrt{S_{xx}} \hat{\beta}_1$

which implies

$$\hat{\beta}_0 = \frac{1}{\sqrt{n}} Z_1 - \frac{\bar{x}}{\sqrt{S_{xx}}} Z_2$$

$$\hat{\beta}_1 = \frac{1}{\sqrt{S_{xx}}} Z_2$$

Since Z_1, Z_2 are independent normal RV, $(\hat{\beta}_0, \hat{\beta}_1) \sim$ bivariate normal.

Besides, we can show that $RSS = \sum_n (Y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \dots$
 $= \sum_{i=3}^n Z_i^2$ (n-2 dof)

$\hat{\sigma}^2 = RSS/n$ is independent of (Z_1, Z_2) , and thus $(\hat{\beta}_0, \hat{\beta}_1)$.

$n \hat{\sigma}^2 / \sigma^2 = RSS / \sigma^2 = \sum Z_i^2 / \sigma^2$, thus it has χ^2 distr w/ n-2 dof.

Distribution of $\hat{\beta}$ (Multiple LR)

Multiple LR, $\hat{\beta}_{d \times 1}$, $X_{n \times d}$ fixed, $y_{n \times 1}$ RV

$$\hat{\beta} = X^+ y, \quad \hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$$

Calculation of $\text{Var}(\hat{\beta})$:

Theorem: If $z = Py$. $\text{Var}(z) = P \text{Var}(y) P^T$

$$\begin{aligned}\text{Proof: } \text{Var}(\hat{\beta}) &= \text{Var}(Py) = E((Py)(Py)^T) - E(Py) E(Py)^T \\ &= P E(yy^T) P^T - P E(y) E(y)^T P^T \\ &= P \text{Var}(y) P^T\end{aligned}$$

Using theorem above,

$$\begin{aligned}\text{Var}(\hat{\beta}) &= (X^T X)^{-1} X^T \cdot \sigma^2 I \cdot X (X^T X)^{-1} \\ &= \sigma^2 I \cdot (X^T X)^{-1}\end{aligned}$$

Joint Distribution of $\hat{\beta}, \hat{\sigma}^2$ (Multiple LR)

- $\hat{\beta} \sim \text{multivariate normal } N(\beta, \sigma^2 (X^T X)^{-1})$
- $\hat{\sigma}^2$ is independent of $\hat{\beta}$
- $n \hat{\sigma}^2 / \sigma^2 \sim \chi^2$ w/ $n-d$ dof.

(Can be proved with similar procedure as in simple LR,
but requires some more advanced technique.)

Testing Hypothesis (Simple LR)

- $H_0: \beta_i = v$ with t-test

WLOG, consider $H_0: \beta_1 = v$:

We know: $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{S_{XX}})$, and we observe $\hat{\beta}_1$.

Then statistic

$$Z = \frac{\hat{\beta}_1 - v}{\sigma / \sqrt{S_{XX}}} \sim N(0, 1).$$

That is enough if true var σ^2 is known.

O/w we need to leverage the fact $n\hat{\sigma}^2 / \sigma^2 \sim \chi^2$.

t-test statistic:

$$Z = \frac{\hat{\beta}_1 - v}{\hat{\sigma}' / \sqrt{S_{XX}}} = \frac{\hat{\beta}_1 - v}{\sqrt{RSS} / \sqrt{(n-2)S_{XX}}}$$

- $H_0: \beta_0 = v_0, \beta_1 = v$, with F-test

First use centered version of LR to get independent $\hat{\beta}_0^*, \hat{\beta}_1^*$.

$$y_i = \beta_0^* + \beta_1^*(x_i - \bar{x}) + \varepsilon_i \quad (\text{In fact, } \beta_1^* = \hat{\beta}_1)$$

$$\hat{\beta}_0^* = \bar{y}, \quad \hat{\beta}_1^* = S_{XY} / S_{XX}$$

$$\hat{\beta}_0^* \sim N(\beta_0^*, \frac{\sigma^2}{N}), \quad \hat{\beta}_1^* \sim N(\beta_1^*, \frac{\sigma^2}{S_{XX}})$$

For known σ^2 , T statistics for $\hat{\beta}_0^*$ and $\hat{\beta}_1^*$ are independent.

Denote them as t_0, t_1 , $t_i \sim \mathcal{N}^2$ w/ 1 dof.

Combine them by $t = t_0 + t_1 \sim \chi^2$ w/ 2 dof.

t contains σ^2 which is unknown. To cancel σ^2 , leverage

$n\hat{\sigma}^2 / \sigma^2 \sim \chi^2_{n-2}$ to get a new statistic

$Z = t / (n\hat{\sigma}^2 / \sigma^2)$ which contains only observable RVs.
 $Z \sim F_{n-2}^2$. (because of the form χ^2 / χ^2)

Testing Hypothesis (Multiple LR)

- $H_0: \beta_i = \nu$

We know that $\hat{\beta}_i \sim N(\beta_i, \text{diag}_i(G^T(X^T X)^{-1}))$

Thus we can construct statistics

$$Z = \frac{\hat{\beta}_i - \nu}{\sigma_i} \sim \chi^2$$

- $H_0: R_{(d-s) \times d} \beta = r$

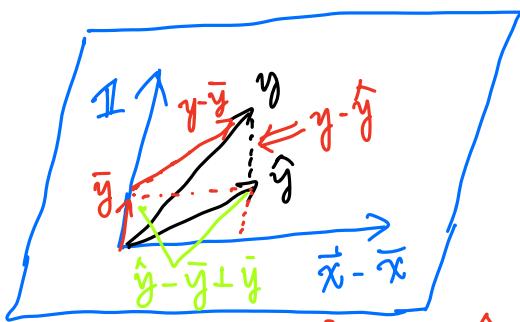
This is a general linear hypothesis that involves multiple dimensions of β . It's possible to construct an F-statistic for this case.

Decomposition of Variation (for simple LR)

Let $SXX = \sum (x_i - \bar{x})(x_i - \bar{x})$, $S\{X,Y\}\{X,Y\}$ can be defined similarly.

$$\text{Now, } SYY = \sum (y_i - \bar{y})^2$$

$$\begin{aligned}
 &= \sum (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\
 &= \underbrace{\sum (y_i - \hat{y}_i)^2}_{\substack{\text{RSS} \\ (\text{residual sum-of-sq})}} + \underbrace{\sum (\hat{y}_i - \bar{y})^2}_{\substack{\text{SSReg} \\ (\text{regression Sum-of-Sq.})}} + 2 \sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \\
 &= \underbrace{\sum (y_i - \hat{y}_i)^2}_{\substack{\text{RSS} \\ (\text{residual sum-of-sq})}} + \underbrace{\sum (\hat{y}_i - \bar{y})^2}_{\substack{\text{SSReg} \\ (\text{regression Sum-of-Sq.})}} + \underbrace{2 \sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})}_{\substack{e_i \\ (\text{error term})}} = \hat{y}_i - \bar{y}_i
 \end{aligned}$$



$$SYY = \|y - \bar{y}\|_2^2 = \|y - \hat{y}\|_2^2 + \|\hat{y} - \bar{y}\|_2^2$$

This is an orthogonal decomposition of SYY .

$\left(\begin{array}{l} \text{also explained} \\ \text{portion of} \\ \text{variability} \\ \text{by regression.} \end{array} \right)$

$$= \langle \vec{e}, \vec{\hat{y}} - \bar{y} \rangle$$

$$= b_1 \langle \vec{e}, \vec{x} - \bar{x} \rangle$$

$$= 0$$

$$\text{By def, } \hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x.$$

$$\text{By choosing } \hat{\beta}_0, \bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}.$$

$$\text{Thus } \hat{y} - \bar{y} = \hat{\beta}_1 (x - \bar{x}).$$

(*)

(will be referred to)

Thus, $SYY = RSS + SSReg$. Note, the decom above is just the orthogonal decomposition with constant dimension removed from every term.

$$SYY = \|y - \bar{y}\|_2^2 \quad \text{fit with constant term} \quad \|y - \bar{y}\|_2^2 - \|\hat{y} - \bar{y}\|_2^2 \quad \text{fit with } x - \bar{x} \text{ term}$$

$$RSS = \|y - \hat{y}\|_2^2 = \|y - \bar{y}\|_2^2 - \|\hat{y} - \bar{y}\|_2^2$$

$$SSReg = \|\hat{y} - \bar{y}\|_2^2 \Rightarrow \text{norm of fitting w/ } x - \bar{x} \text{ term}$$

Analysis of Variance (ANOVA)

(used to compare the means of multiple normal distr)
(also useful for testing adequacy of LR)

Setting (One-way layout)

Random samples of d normal distr are available.

Variance of each distr is σ^2 , we want to compare their means given observed values of the samples.

Notation: For $i=1 \dots d$, we observe $Y_{i1}, Y_{i2} \dots Y_{in_i} \sim N(\mu_i, \sigma^2)$

To find μ_i , we can solve the following:

$$Y = \begin{bmatrix} Y_{11} \\ \vdots \\ Y_{1n_1} \\ Y_{21} \\ \vdots \\ Y_{2n_2} \\ \vdots \\ Y_{dn_d} \end{bmatrix} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & & \vdots & & \vdots \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix}$$

which is a multiple regression.

Let \bar{Y}_{it} be sample mean of n_j i -th population,

$$\text{i.e., } \bar{Y}_{it} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}.$$

$$\text{MLE of } G^2, \hat{G}^2 = \frac{1}{n} \sum_{i=1}^d \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{it})^2$$

(sample var using mean of each corresponding population)

Target Hypothesis

$$H_0: \mu_1 = \mu_2 = \dots = \mu_d$$

Partitioning of Sum of Squares

$$\text{Let } \bar{Y}_{++} = \frac{1}{n} \sum \sum Y_{ij}, \quad S_{\text{tot}}^2 = \sum_{i=1}^d \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{++})^2$$

We can partition S_{tot} as:

$$S_{\text{tot}}^2 = S_{\text{Res}}^2 + S_{\text{Betw}}^2$$

where

$$S_{\text{Res}}^2 = \sum \sum (Y_{ij} - \bar{Y}_{ii})^2$$

$$S_{\text{Betw}}^2 = \sum_{i=1}^d n_i (\bar{Y}_{ii} - \bar{Y}_{++})^2$$

variation that can't be explained by its population identity

(in-population variation)

(between-population variation)

variation explained by population id

Besides, $S_{\text{Res}}^2 / \sigma^2 \sim \chi^2$ w/ $n-d$ dof,

S_{Res} and S_{Betw} are independent.

Also, if $\mu_1 = \dots = \mu_d$, which is a useful scenario to test,

$S_{\text{Betw}}^2 / \sigma^2 \sim \chi^2$ w/ $d-1$ dof.

Intuition of $d-1$ dof: S_{Betw}^2 is of form $\sum (x_i - \bar{x})^2$
can use similar argument when proving sample var $\frac{n \sigma^2}{\sigma^2} \sim \chi^2_{n-1}$.

(ANOVA table)

This partitioning can be summarized as a table:

Source of Variation	DOF	Sum of Square	Mean Square
Between population	$d-1$	S_{Betw}^2	$S_{\text{Betw}}^2 / (d-1)$
Residual	$n-d$	S_{Res}^2	$S_{\text{Res}}^2 / (n-d)$
Total	$n-1$	S_{tot}^2	

Mean Square: Sum of square / DOF, used in hypothesis testing.

Note: Residual mean square is σ^2 of regression.

ANOVA Test

Let test statistic

$$U^2 = \frac{S_{\text{Betw}}^2 / (d-1)}{S_{\text{Res}}^2 / (n-d)}, \quad U^2 \sim F(d-1, n-d).$$

Intuition: If $\mu_i \neq \mu_j$, S_{Betw}^2 will become large, and U^2 fall into reject region.

ANOVA and LR

As developed earlier, ANOVA is just a special type of multiple LR, with design matrix of specific form.

The partition of variation in ANOVA is backed by that of LR.

Thus the F-test statistics also has its counterpart in LR.

(checking adequacy of LR)

ANOVA as Significance Test in LR

When we apply LR, we made the linear assumption which is not warranted.

We can use ANOVA as a significance test for LR.

$$S_{\text{tot}}^2 \leftrightarrow SYY \quad S_{\text{Betw}}^2 \leftrightarrow SS_{\text{Reg}} \quad S_{\text{Res}}^2 \leftrightarrow RSS$$

Thus we can develop ANOVA style test for LR:

$$H_0: \beta \text{ entries except for the constant term equals } 0. \quad (\# \text{ group} = 1 \text{ in ANOVA})$$
$$U^2 = \frac{SS_{\text{Reg}} / \text{dof}(SS_{\text{Reg}})}{RSS / \text{dof}(RSS)} = \frac{SS_{\text{Reg}} / (d-1)}{RSS / (n-d)}$$

Intuition: SS_{Reg} / RSS measures how much of variation is explained by LR and thus reflects significance of regression.

R^2 (coefficient of determination) and Goodness of Fit (Simple LR)

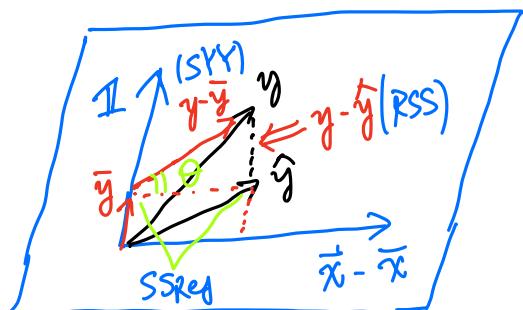
$$R^2 = 1 - \frac{RSS}{SYY} = \frac{SS_{Reg}}{SYY}$$

ratio of unexplained var

$$R^2 \in [0, 1]$$

measures goodness of fit.

LR w/ intercept



$$SYY = \|y - \bar{y}\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - \bar{y}\|^2$$

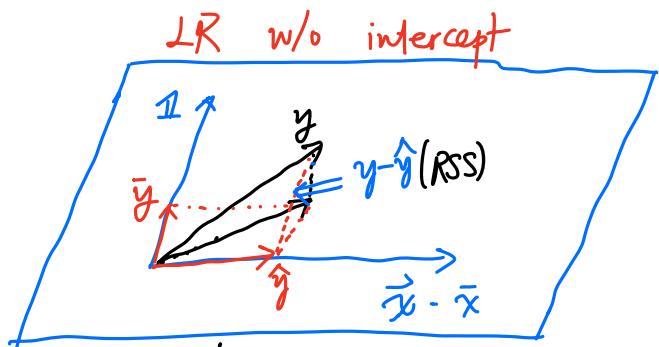
$$R^2 = 1 - \frac{RSS}{SYY} = \frac{SS_{Reg}}{SYY}$$

This step, as derived above (*), requires being able to choose \hat{y} such that $\langle \vec{e}, \hat{y} - \bar{y} \rangle = 0$.

Thus $R^2 = \frac{SS_{Reg}}{SYY}$ is only valid for LR with nonzero intercept.

For LR with no intercept,

$$R^2 = 1 - \frac{RSS}{\sum y_i^2}$$



$$R^2 = \frac{\|\hat{y}\|}{\|y\|} = 1 - \frac{RSS}{\sum y_i^2}$$

Goodness of Fit for Multiple Regression

$$R^2 = 1 - \frac{RSS}{SYY}$$

$$\begin{aligned} \text{where } RSS &= \|y - \hat{y}\|_2^2 \\ &= y^T y - y^T (X \hat{\beta}) \\ &= y^T y - y^T \boxed{X X^T} y \\ &= SYY - SS_{Reg} \end{aligned}$$

$$SYY = \|y - \bar{y}\|_2^2$$

orthogonal projection of y on $\text{span}(X)$

$y^T \cdot \text{proj}(y) \Rightarrow$ component of $\|y\|$ -norm
that's covered by $\text{span}(X)$

\Rightarrow amount explained by reg
 $\Rightarrow SS_{Reg}$

(Decomposition of err arise naturally from geometric perspective)

Adjusted \bar{R}^2 : $\bar{\bar{R}}^2$

$$\bar{\bar{R}}^2 = 1 - \frac{RSS/(n-d)}{SYY/(n-1)} = 1 - \frac{n-1}{n-d} (1 - \bar{R}^2)$$

obtained from dividing RSS and SYY with their degree of freedom.

R^2 and Pearson Correlation (For simple LR)

Note, $R^2 = \text{Corr}(\vec{x}, \vec{y})$.

$$\begin{aligned}\text{Proof: } R^2 &= SS_{\text{Reg}} / SYY \\ &= \sum (\hat{y}_i - \bar{y})^2 / \sum (y_i - \bar{y})^2 \\ &= \sum (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{y})^2 / \sum (y_i - \bar{y})^2 \\ &= \sum ((\bar{y} - \hat{\beta}_1 \bar{x}) + \hat{\beta}_1 x_i - \bar{y})^2 / \sum (y_i - \bar{y})^2 \\ &= \hat{\beta}_1^2 \sum (x_i - \bar{x})^2 / \sum (y_i - \bar{y})^2 \\ &= \frac{SXY^2}{SXX^2} \cdot SXX / SYY \\ &= SXY^2 / (SXX \cdot SYY) \\ &= \left[SXY / (\sqrt{SXX} \cdot \sqrt{SYY}) \right]^2 \\ &= \text{Corr}^2(x, y)\end{aligned}$$

Alt. Definition of R^2 (For simple LR)

R^2 is correlation of \vec{x} and \vec{y} squared. $R^2 = \text{corr}^2(\vec{x}, \vec{y})$

For general LR, $R^2 = \text{Corr}^2(y, \hat{y})$.

Multicollinearity

Multiple Regression from Simple Reg

Let \tilde{x}_{i,nx_1} ($i=1 \dots d$) be i -th dim vector, if $\langle x_i, x_j \rangle = 0$,

$$\begin{aligned}\hat{\beta} &= (\tilde{x}^T \tilde{x})^{-1} \tilde{x}^T y = \left(\begin{pmatrix} x_1^T \\ \vdots \\ x_d^T \end{pmatrix} (x_1 \dots x_d) \right)^{-1} \begin{pmatrix} x_1^T y \\ \vdots \\ x_d^T y \end{pmatrix} \\ &= \begin{pmatrix} \frac{1}{x_1^T x_1} & & & \\ & \ddots & & 0 \\ & & \ddots & \dots \\ 0 & & & \frac{1}{x_d^T x_d} \end{pmatrix} \cdot \begin{pmatrix} x_1^T y \\ \vdots \\ x_d^T y \end{pmatrix}\end{aligned}$$

Or, $\hat{\beta}_i = \frac{\langle x_i, y \rangle}{\langle x_i, x_i \rangle}$ same as simple reg of y on x_i .

i.e., orthogonal $x_i \Rightarrow$ can reduce multiple reg to simple reg on each dimension.

For general case (non orthogonal), X_i 's can be orthogonalized, e.g., w/ Gram-Schmidt. upper triangular

Let $X_{n \times d} = Q_{n \times d} R_{d \times d}$ $\leftarrow \text{diag}(R) = I$ for Gram-Schmidt

$Q = (z_1, \dots, z_d)$, where z_i is residual of regressing X_i on $X_{<i}$.

Note, if consider constant term, set $X_0 = z_0 = I$.

$$\text{Now, } \hat{\beta} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \bar{y} = \bar{R}^{-1} \bar{Q}^T \bar{y}$$

$$\hat{y} = \bar{Q} \bar{Q}^T \bar{y}$$

$$\text{Consider last dim } \hat{\beta}_d, \quad \hat{\beta}_d = \frac{\langle z_d, y \rangle}{\langle z_d, z_d \rangle}$$

We consider $\hat{\beta}_d$ here,
but no loss of generality
because we can rearrange
to let any X_i be X_d .

$$\begin{aligned} \text{Var}(\hat{\beta}_d) &= E(\hat{\beta}_d^2) - E(\hat{\beta}_d)^2 \\ &= \frac{E(z_d^T y y^T z_d) - [z_d^T E(y)]^2}{(z_d^T z_d)^2} = \frac{z_d^T [E(y y^T) - E(y) E(y)^T] z_d}{(z_d^T z_d)^2} \\ &= \frac{z_d^T \text{cov}(y, y) z_d}{(z_d^T z_d)^2} = \frac{\sigma^2}{\|z_d\|_2^2} \end{aligned}$$

Now, if z_d is correlated with some X_i , $\|z_d\|$ will be small and $\text{Var}(\hat{\beta}_d)$ will be large.

Dr, more generally, the precision with which we can estimate $\hat{\beta}_i$ depends on how much X_i is unexplained by other X_j 's.

Multicollinearity

Definition: X being rank deficient, or feature vectors are not linearly independent. This will cause LR estimation to have large variance.

The Canonical Form

(useful when X is rank deficient, causing $X^T X$ to be singular)

Consider spectral decomposition of symmetric matrix

$$X^T X = P \Lambda P^T \quad \begin{array}{l} \text{(symmetric matrix can be orthogonally} \\ \text{diagonalized)} \end{array}$$
$$P = (p_1, \dots, p_d), \quad P P^T = I$$

$$\text{The model becomes: } y = X P P^T \beta + \varepsilon$$
$$= \tilde{X} \tilde{\beta} + \varepsilon$$

where $\tilde{X} = X P$, $\tilde{\beta} = P^T \beta$, and $\tilde{X}^T \tilde{X} = P^T X^T X P = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$,
which means columns of \tilde{X} are orthogonal.

Note: not $\hat{\beta}$

$$\text{Var}(\hat{\beta}) = \sigma^2 (X^T X)^{-1} = \sigma^2 \sum \lambda_i^{-1} p_i p_i^T$$

$$\hat{\beta} = X^+ y = \Lambda^{-1} \tilde{X}^T \tilde{y}$$

$$\text{Var}(\hat{\beta}) = \sigma^2 \Lambda^{-1}$$

Handling Exact Multicollinearity

Consider canonical form $y = \tilde{X}\tilde{\beta} + \varepsilon$

$\tilde{x}_i = Xp_i$ is the i -th principle component

$$\tilde{x}_i^T \tilde{x}_i = p_i^T X^T X p_i = \lambda_i.$$

Assume $\text{rank}(X) = d-k$, then $\lambda_{d-k+1} \dots \lambda_d = 0$.

Divide the matrices by $d-k, k$:

$$P = (P_1, P_2) \quad A = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \quad \tilde{X} = (\tilde{X}_1, \tilde{X}_2) = (XP_1, XP_2)$$

$$\tilde{\beta} = \begin{pmatrix} \tilde{\beta}_1 \\ \tilde{\beta}_2 \end{pmatrix} = \begin{pmatrix} P_1^T \beta \\ P_2^T \beta \end{pmatrix}, \quad \tilde{x}_2 = 0 \text{ because } \tilde{x}_i^T \tilde{x}_i = 0 \text{ for } \lambda_{d-k+1} \dots \lambda_d.$$

$$\begin{aligned} y &= \tilde{x}_1 \tilde{\beta}_1 + \tilde{x}_2 \tilde{\beta}_2 + \varepsilon \\ &= \tilde{x}_1 \tilde{\beta}_1 \end{aligned}$$

$$\hat{\tilde{\beta}}_1 = \tilde{X}_1^+ y, \quad \hat{\tilde{\beta}} = \begin{pmatrix} \hat{\tilde{\beta}}_1 \\ 0 \end{pmatrix} = \hat{X}^+ y = (P \Lambda^{-1} P^T) X^T y$$

where $\Lambda^{-1} = \begin{pmatrix} \Lambda_1^{-1} & 0 \\ 0 & 0 \end{pmatrix}$

Soft Solution to Multicollinearity

Includes Ridge, Shrinkage, Lasso ...

Shrinkage Methods

Ridge (L2) Regression

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \sum (y_i - \beta_0 - \beta^T x_i)^2 + \lambda \|\beta\|_2^2 \right\}$$

β do not have 0 index

Dual form:
$$\begin{aligned} & \underset{\beta}{\operatorname{argmin}} \sum (y_i - f_0 - \beta^T x_i)^2 \\ & \text{s.t. } \|\beta\|_2^2 \leq t \end{aligned}$$

where there is 1 to 1 correspondence between λ and t .

Why regularization:

When there are correlated variables, there may exist large coeff for one var cancelled by large negative coeff of another variable, causing large variance.

Note, ridge is not equivalent under scaling of inputs.
Thus need to standardize before fitting.

Solution of Ridge

$$\text{From primal form, } L(\beta) = (y - X\beta)^T(y - X\beta) + \lambda \beta^T \beta \\ \Rightarrow \hat{\beta}^* = (X^T X + \lambda I)^{-1} X^T y$$

Interpretation from SVD:

$$\text{Let } X = UDV^T,$$

$$\hat{\beta}_{OLS} = (X^T X)^{-1} X^T y = UU^T y = \sum u_j u_j^T y$$

$$\hat{\beta}_{RIDGE} = X\beta^* = \sum_{j=1}^d u_j \frac{d_j}{d_j^2 + \lambda} u_j^T y$$

shrink SVD coordinate by $\frac{d_j^2}{d_j^2 + \lambda}$
dimension w/ smaller d_j is shrinked more.

Lasso (L1)

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \left\{ \frac{1}{2} \sum (y_i - \beta_0 - \beta^T x_i)^2 + \lambda \|\beta\|_1 \right\}$$

Dual:

$$\underset{\beta}{\operatorname{argmin}} \quad \sum (y_i - \beta_0 - \beta^T x_i)^2 \\ \text{s.t.} \quad \|\beta\|_1 \leq t$$

Making t sufficiently small will cause certain dimension of β to be exactly zero.

Varying $t \Rightarrow$ continuous subset selection

Subset Selection vs. Shrinkage Method

For orthogonal X , $\hat{\beta}_j = :$

Best subset (size M)

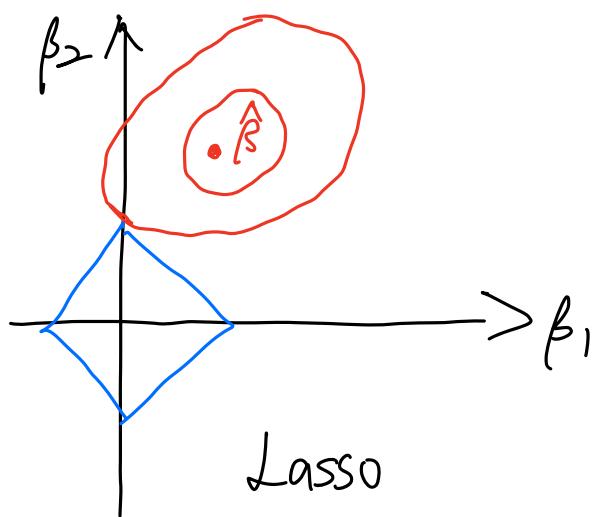
Ridge

Lasso

$$\hat{\beta}_j \cdot \begin{cases} 1 & \|\hat{\beta}_j\| \geq \|\hat{\beta}_M\| \\ \hat{\beta}_j / (1+\lambda) & \|\hat{\beta}_j\| < \|\hat{\beta}_M\| \end{cases}$$

$$\text{sign}(\hat{\beta}_j) \cdot (|\hat{\beta}_j| - \lambda)_+$$

Illustration:



— : constraint

— : loss function contour

Note, the illustration shows intuitively why L_1 (Lasso) tend to completely eliminate smaller dimensions (due to shape of contour)