

Bayesian inference in simple conjugate families.

Preview: All you need to know

- ① Bayes' rule: $p(\theta|x) = \frac{P(x|\theta)P(\theta)}{P(x)} \propto P(x|\theta)P(\theta)$ referred to as Bayes' rule*
- ② Gaussian kernel: $\exp\{-\frac{w}{2}(x-\mu)^2\}$
- ③ Beta kernel: $x^{\alpha-1}(1-x)^{\beta-1}$
- ④ Gamma kernel: $x^{\alpha-1}\exp\{-\beta x\}$
- ⑤ Two-step method for finding posteria:
 - A. Write the kernel of $p(\theta|x)$ using tip 1;
 - B. match this kernel with an exist distribution to derive normalization constant.

(A) $p(w|x_1, \dots, x_N) \propto P(x_1, \dots, x_N|w) \cdot P(w)$ (tip 1)

Beta-Bernoulli $= \left[\prod_{i=1}^N P(x_i|w) \right] P(w)$ (x_1, \dots, x_N ind.)

$$\propto \left[\prod_{i=1}^N w^{x_i} (1-w)^{(1-x_i)} \right] w^{a-1} (1-w)^{b-1}$$

$$= w^{\sum_{i=1}^N x_i + a - 1} (1-w)^{N - \sum_{i=1}^N x_i + b - 1}$$

which is the kernel of Beta($\sum_{i=1}^N x_i + a$, $N + b - \sum_{i=1}^N x_i$)

(B) $\begin{cases} y_1 = x_1 / (x_1 + x_2) \\ y_2 = x_1 + x_2 \end{cases} \Rightarrow \begin{cases} x_1 = y_1 y_2 \\ x_2 = y_2 - y_1 y_2 \end{cases}$

$$J = \begin{vmatrix} \partial x_1 / \partial y_1 & \partial x_1 / \partial y_2 \\ \partial x_2 / \partial y_1 & \partial x_2 / \partial y_2 \end{vmatrix} = \begin{vmatrix} y_2 & y_1 \\ -y_2 & 1 - y_1 \end{vmatrix} = y_2$$

Tip for Jacobian: consider the integrators are arranged in a column vector.
originally, we're integrating on $(dx_1, dx_2)'$, since

$$\begin{bmatrix} dx_1 \\ dx_2 \end{bmatrix} = \begin{bmatrix} \partial x_1 / \partial y_1 & \partial x_1 / \partial y_2 \\ \partial x_2 / \partial y_1 & \partial x_2 / \partial y_2 \end{bmatrix} \begin{bmatrix} dy_1 \\ dy_2 \end{bmatrix}$$

and the new integration is taken on $(dy_1, dy_2)'$, the Jacobian is sort of "amount of scaling", which is the determinant of the transformation mat

Back to problem (B). With Jacobian J , we have

$$f_Y(y_1, y_2) = f_X(x_1, x_2, y_2 - x_1 y_2) \cdot |J|,$$

Calculations omitted, finally, we have

$$f_Y(x_1) = \int f_Y(x_1, x_2) dx_2 \sim \text{Ga}(a_1 + a_2, 1)$$

$$f_Y(x_2) = \int f_Y(x_1, x_2) dx_1 \sim \text{Beta}(a_1, a_2)$$

So we learn from this problem:

① $\text{Ga}(a_1, 1) / [\text{Ga}(a_1, 1) + \text{Ga}(a_2, 1)] \sim \text{Ga}(a_1 + a_2, 1)$

② Gamma \mapsto Beta: $\text{Ga}(a_1, 1) + \text{Ga}(a_2, 1) \sim \text{Beta}(a_1, a_2)$

(c) Normal - Normal

$$p(\theta) \propto \exp \left\{ -\frac{(\theta - m)^2}{2\nu} \right\},$$

$$\begin{aligned} p(x|\theta) &\propto \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2 \right\} = \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n [(x_i - \bar{x})^2 + (\bar{x} - \theta)^2] \right\} \\ &= \exp \left\{ -\frac{1}{2\sigma^2} [S_x + n(\bar{x} - \theta)^2] \right\}. \end{aligned}$$

Here \underline{x} denotes vector $(x_1, x_2, \dots, x_n)^T$, $S_x = \sum_{i=1}^n (x_i - \bar{x})^2$ and \bar{x} is the mean value of elements of \underline{x} . By Bayes' rule*,

$$\begin{aligned} p(\theta|\underline{x}) &\propto \exp \left\{ -\frac{(\theta - m)^2}{2\nu} - \frac{S_x}{2\sigma^2} - \frac{n(\bar{x} - \theta)^2}{2\sigma^2} \right\} \\ &= \exp \left\{ \theta^2 \left(-\frac{1}{2\nu} - \frac{n}{2\sigma^2} \right) + \theta \left(\frac{m}{\nu} + \frac{n\bar{x}}{\sigma^2} \right) \right\} \\ &\propto \exp \left\{ -\frac{1}{2} \cdot \left[\frac{1}{\nu} + \frac{n}{\sigma^2} \right] \cdot \left(\theta - \frac{m/\nu + n\bar{x}/\sigma^2}{1/\nu + n/\sigma^2} \right)^2 \right\} \\ &\sim \mathcal{N} \left(\frac{m/\nu + n\bar{x}/\sigma^2}{1/\nu + n/\sigma^2}, \text{precision} = \frac{1}{\nu} + \frac{n}{\sigma^2} \right) \end{aligned}$$

The normal posterior is a combination of prior and likelihood, with 1. additive precision, i.e., precision of posterior equals the sum of precisions of prior and likelihood; 2. convex combined mean, i.e., the mean of output distribution is a convex combination of the means of component distributions, with each weight proportional to precision.

(D) Inverse Gamma - Normal.

$$p(w) \propto w^{a-1} e^{-bw}$$

$$p(x|w) \propto \exp \left\{ -\frac{w}{2} [S_x + n(\bar{x} - \theta)^2] \right\} \cdot w^{n/2}$$

$$\Rightarrow p(w|x) \propto w^{\frac{n}{2}+a-1} \exp \left\{ -\left[\frac{S_x + n(\bar{x} - \theta)^2}{2} + b \right] w \right\}$$

$$\sim \text{Gamma} \left(\frac{n}{2} + a, \frac{S_x + n(\bar{x} - \theta)^2}{2} + b \right)$$

(E) Normal - Normal, generalized sample sd.

Following steps in problem (c), we have

$$p(\theta|x) \propto \exp \left\{ -\frac{1}{2} \left[\frac{1}{v} + \sum_{i=1}^n \frac{1}{\sigma_i^2} \right] \cdot \left(\theta - \frac{m/v + \sum_{i=1}^n x_i/\sigma_i^2}{1/v + \sum_{i=1}^n 1/\sigma_i^2} \right)^2 \right\}$$

$$\sim N \left(\frac{m/v + \sum_{i=1}^n x_i/\sigma_i^2}{1/v + \sum_{i=1}^n 1/\sigma_i^2}, \text{precision} = \frac{1}{v} + \sum_{i=1}^n \frac{1}{\sigma_i^2} \right).$$

The result gives a more clear demonstration of the conclusion we reached at (c).

(F) $p(x, w) = p(x|w) p(w)$

$$\propto w^{1/2} \exp \left\{ -\frac{w}{2} (x-m)^2 \right\} w^{a/2-1} \exp \left\{ -\frac{b}{2} w \right\}$$

$$= w^{\frac{a+1}{2}-1} \exp \left\{ -\frac{(x-m)^2 + b}{2} w \right\},$$

which is a Gamma function w.r.t. w , so when we integral out w , what is left will be

$$p(x) = \int_0^{+\infty} p(x, w) dw \propto \left[\frac{b}{2} + \frac{(x-m)^2}{2} \right]^{-\frac{a+1}{2}}$$

$$\propto \left[1 + \frac{1}{a} \cdot \frac{(x-m)^2}{b/a} \right]^{-\frac{a+1}{2}}, \text{ which is the kernel of } t \text{ distribution,}$$

i.e., $v = a$, $m = m$, $s^2 = b/a$ ← scale

↑ center

degrees of freedom

t distribution as a scale mixture of Gaussian:

$$p(x) = \int p(x|w) p(w) dw \hat{=} \sum \underbrace{Pr(x|w \in \Delta w_i)}_{\text{mixture component}} \cdot \underbrace{Pr\{w \in \Delta w_i\}}_{\text{weight}}.$$

We use this property to deal with outliers.

The multivariate normal distribution. - Basics

$$(A) \text{Cov}(\underline{x}) = E(\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})' = E(\underline{x}\underline{x}') - E(\underline{x})\underline{\mu}' - \underline{\mu}E(\underline{x})' + \underline{\mu}\underline{\mu}' \\ = E(\underline{x}\underline{x}') - 2\underline{\mu}\underline{\mu}' + \underline{\mu}\underline{\mu}' = E(\underline{x}\underline{x}') - \underline{\mu}\underline{\mu}'$$

$$(B) f(\underline{z}) = f(z_1)f(z_2)\dots f(z_p) \quad (\text{independence of } z_1 \sim z_p) \\ = (2\pi)^{-p/2} \exp\left\{-\frac{1}{2}\sum_{i=1}^p z_i^2\right\} = (2\pi)^{-p/2} \exp\left\{-\frac{1}{2}\underline{z}'\underline{z}\right\} \quad (\text{PDF})$$

$$M_{\underline{z}}(\underline{t}) = E[\exp\{\underline{t}'\underline{z}\}] = E[\exp\{t_1 z_1\} \exp\{t_2 z_2\} \dots \exp\{t_p z_p\}] \\ = E[\exp\{t_1 z_1\}] \cdot E[\exp\{t_2 z_2\}] \dots E[\exp\{t_p z_p\}] \quad (\text{MGF})$$

$$(C) \Rightarrow : \underline{x} \sim N(\underline{\mu}, \underline{\Sigma}), \text{ then } \forall \underline{a} \in \mathbb{R}^p / \{0\}, \underline{a}'\underline{x} \text{ is univariate normal,} \\ \text{i.e., MGF of } z = \underline{a}'\underline{x} \text{ is } M_z(t) = \exp\left\{E(z)t + \frac{1}{2}\text{var}(z)t^2\right\}, \text{ where} \\ E(z) = E(\underline{a}'\underline{x}) = \underline{a}'\underline{\mu}, \text{ var}(z) = \underline{a}'\underline{\Sigma}\underline{a}. \text{ so} \\ M_z(t) = E(\exp\{zt\}) = \exp\left\{\underline{a}'\underline{\mu}t + \frac{1}{2}\underline{a}'\underline{\Sigma}\underline{a}t^2\right\}.$$

On the other hand,

$$M_{\underline{z}}(\underline{\lambda}) = E[\exp\{(\underline{\lambda})'\underline{z}\}] = \exp\left\{(\underline{\lambda})'\underline{\mu} + \frac{1}{2}(\underline{\lambda})'\underline{\Sigma}(\underline{\lambda})\right\} \\ = \exp\left\{\underline{t}'\underline{\mu} + \frac{1}{2}\underline{t}'\underline{\Sigma}\underline{t}\right\}, \quad \underline{t} = \underline{\lambda},$$

which means

$$M_{\underline{x}}(\underline{t}) = E[\exp\{\underline{t}'\underline{x}\}] = \exp\left\{\underline{t}'\underline{\mu} + \frac{1}{2}\underline{t}'\underline{\Sigma}\underline{t}\right\}$$

$$\Leftarrow : \text{If } \underline{x} \text{ has MGF } M_{\underline{x}}(\underline{t}) = \exp\left\{\underline{t}'\underline{\mu} + \frac{1}{2}\underline{t}'\underline{\Sigma}\underline{t}\right\}, \text{ then for } \forall \underline{a} \in \mathbb{R}^p / \{0\} \\ \text{and } z = \underline{a}'\underline{x},$$

$$M_z(\lambda) = E[\exp\{\lambda z\}] = E[\exp\{(\lambda \underline{a})'\underline{x}\}] \\ = M_{\underline{x}}(\underline{t} = \lambda \underline{a}) = \exp\left\{\lambda \underline{a}'\underline{\mu} + \frac{1}{2}\lambda^2 \underline{a}'\underline{\Sigma}\underline{a}\right\}$$

Since z has mean $\underline{a}'\underline{\mu}$ and variance $\underline{a}'\underline{\Sigma}\underline{a}$, the MGF of z can be rewritten by $M_z(t) = \exp\left\{\lambda E(z) + \frac{1}{2}\lambda^2 \text{var}(z)\right\}$, which indicates z is a univariate Gaussian variable.

Above proves $\left\{ \begin{array}{l} \text{"any linear combination of } \underline{x}'\text{s components is univar normal"} \\ = \text{"} \underline{x} \text{ is multivar normal"} \end{array} \right\}$

$$\Leftrightarrow \left\{ \begin{array}{l} \text{"} \underline{x} \sim N(\underline{\mu}, \underline{\Sigma}) \text{"} \\ = \text{"MGF of } \underline{x} \text{ is } \exp\left\{\underline{t}'\underline{\mu} + \frac{1}{2}\underline{t}'\underline{\Sigma}\underline{t}\right\}" \end{array} \right\}$$

(D) I'm skeptical about the "full col rank" statement here, in that as far as I can see, it is more appropriate for L to be "full row rank" than "full col rank".

Basically, the idea is to show that $\forall \underline{a} \in \mathbb{R}^p \setminus \{0\}$, $\underline{b}' = \underline{a}'L$ is also non-zero, which requires the ROWS of L being linearly independent.

(E) Without loss of generality, let \underline{X} follows a non-singular normal distribution, i.e., $\underline{X} \sim \mathcal{N}(\underline{\mu}, \Sigma)$ and Σ is non-singular.
 $\Rightarrow \Sigma$ is positive definite. Then we can apply spectral decomposition to Σ . Let $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_p$ are eigen vectors of Σ , then:

i) Σ is symmetric \Rightarrow eigen vectors with different eigen values are orthogonal;

ii) for eigen vectors with the same eigen value, apply Gram-Schmidt orthogonalization;

iii) normalize all the eigen vectors.

At this stage we will have an orthogonal mat $P = [\underline{v}_1, \underline{v}_2, \dots, \underline{v}_p]$, where $\|\underline{v}_i\| = 1$, $\underline{v}_i' \underline{v}_j = 0$, $\forall i \neq j \in \{1, 2, \dots, p\}$, and

$$\Sigma = P \Lambda P' \quad (\text{spectral decomposition}),$$

$$\Lambda = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_p\}.$$

Consider the affine transformation

$\underline{X}^* = L\underline{Z} + \underline{\mu}^*$, where \underline{Z} is standard multivariate normal, $\underline{\mu}^* = \underline{\mu}$ and $L = P\Lambda^{1/2}$, then by what we proved in (D),

$$\text{rank}(P\Lambda^{1/2}) = p \Rightarrow \underline{X}^* \sim \mathcal{N}(\underline{\mu}^*, (P\Lambda^{1/2})(P\Lambda^{1/2})') = \mathcal{N}(\underline{\mu}^*, P\Lambda P').$$

As mean and var characterize a normal distribution, so \underline{X}^* is just \underline{X} .

In this way, for any $\underline{X} \sim \mathcal{N}(\underline{\mu}, \Sigma)$, we can construct the affine-transformation $(L, \underline{\mu})$ s.t. $\underline{X} = L\underline{Z} + \underline{\mu}$.

(F) For $\underline{x} \sim N(\underline{\mu}, \Sigma)$, we know from (E) that we are able to transform them back to standard normal variables:

$$\underline{x} = L\underline{z} + \underline{\mu} \Leftrightarrow \underline{z} = L^{-1}(\underline{x} - \underline{\mu}), \quad L^{-1} = \Lambda^{-1/2} P'$$

Using the R.V. transformation property,

$$f_{\underline{x}}(\underline{x}) = f_{\underline{z}}(\underline{z}(\underline{x})) \cdot |J|_{\underline{x} \rightarrow \underline{z}}$$

$$= (2\pi)^{-p/2} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu})' P \Lambda^{1/2} \Lambda^{1/2} P' (\underline{x} - \underline{\mu}) \right\} \cdot |\Lambda^{-1/2} P'|$$

$$= (2\pi)^{-p/2} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu})' \Sigma^{-1} (\underline{x} - \underline{\mu}) \right\} |\Lambda^{-1/2}| \cdot |P|$$

$$= \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\underline{x} - \underline{\mu})' \Sigma^{-1} (\underline{x} - \underline{\mu}) \right\} \quad (*)$$

$$(*) : |P| = 1, |\Sigma| = |P \Lambda P'| = |P| \cdot |\Lambda| \cdot |P'| = |\Lambda|, |\Lambda^{1/2}| = |\Lambda|^{1/2}$$

$$\Rightarrow |\Lambda^{-1/2}| = |\Lambda|^{-1/2} = |\Sigma|^{-1/2}$$

(G) Both matrices A and B are full ROW rank, not full col rank.

Let the row number of both A and B be p , then $\text{rank}(A) = \text{rank}(B) = p$.

A, B full row rank \Leftrightarrow Any $\underline{L} \in \mathbb{R}^p / \{0\}$, $\underline{L}'A$ and $\underline{L}'B$ are not zero.

By the definition of "multivariate normal" in (c), it's safe to claim

$\underline{x}_1, \underline{x}_2$ are multivariate Gaussian } Any $\underline{L} \in \mathbb{R}^p / \{0\}$, $(\underline{L}'A)\underline{x}_1$ and $(\underline{L}'B)\underline{x}_2$ are univariate normal R.V.'s
 A, B are of full row rank

\Leftrightarrow both $A\underline{x}_1$ and $B\underline{x}_2$ are multivariate normal.

\Leftrightarrow MGF of $A\underline{x}_1$ is $\exp \left\{ \underline{t}'(A\underline{\mu}_1), \frac{1}{2} \underline{t}'(A\underline{\Sigma}_1 A') \underline{t} \right\}$

MGF of $B\underline{x}_2$ is $\exp \left\{ \underline{t}'(B\underline{\mu}_2), \frac{1}{2} \underline{t}'(B\underline{\Sigma}_2 B') \underline{t} \right\}$

Beyond that, \underline{x}_1 is independent with \underline{x}_2 , indicating that $A\underline{x}_1$ is independent with $B\underline{x}_2$, then MGF of $A\underline{x}_1 + B\underline{x}_2$ is:

$$M_{A\underline{x}_1 + B\underline{x}_2}(\underline{t}) = M_{A\underline{x}_1}(\underline{t}) \cdot M_{B\underline{x}_2}(\underline{t}) = \exp \left\{ \underline{t}'(A\underline{\mu}_1 + B\underline{\mu}_2) + \frac{1}{2} \underline{t}'(A\underline{\Sigma}_1 A' + B\underline{\Sigma}_2 B') \underline{t} \right\}$$

From the MGF of $A\underline{x}_1 + B\underline{x}_2$, we know it's a multivariate Gaussian R.V. with mean $A\underline{\mu}_1 + B\underline{\mu}_2$ and variance $A\underline{\Sigma}_1 A' + B\underline{\Sigma}_2 B'$.

The multivariate normal distribution - Conditionals and marginals.

- (A) "Any marginal distribution of a partition of \underline{x} is normal distribution if $\underline{x} \sim N(\underline{\mu}, \Sigma)$ ".

Let $A = [I_k \mid 0_{k \times (p-k)}]$. Clearly A has size $k \times p$ and is of full row rank. By what we have proven in (G) of last section, we have

$$\underline{x}_1 = A\underline{x} \text{ follows } N(A\underline{\mu}, A\Sigma A') = N(\underline{\mu}_1, \Sigma_{11})$$

- (B) As long as \underline{x} follows a non-singular multivariate normal distribution, it's safe to assert that both Σ_{11} and Σ_{22} are invertible.

Then by the inverse of block matrix:

$$\begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}^{-1} = \begin{bmatrix} (\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})^{-1} & -(\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})^{-1}\Sigma_{12}\Sigma_{22}^{-1} \\ -(\Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})^{-1}\Sigma_{21}\Sigma_{11}^{-1} & (\Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})^{-1} \end{bmatrix}$$

we have

$$\Omega_{11} = (\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})^{-1}, \quad \Omega_{12} = -\Omega_{11}\Sigma_{12}\Sigma_{22}^{-1}$$

$$\Omega_{22} = (\Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12})^{-1}, \quad \Omega_{21} = -\Omega_{22}\Sigma_{21}\Sigma_{11}^{-1}$$

$$(c) \log f(\underline{x}_1 | \underline{x}_2) = \log f(\underline{x}_1, \underline{x}_2) - \log f(\underline{x}_2)$$

$$= c + \left(-\frac{1}{2}(\underline{x} - \underline{\mu})' \Omega (\underline{x} - \underline{\mu})\right) - \left(-\frac{1}{2}(\underline{x}_2 - \underline{\mu}_2)' \Omega_{22}(\underline{x}_2 - \underline{\mu}_2)\right)$$

$$= c - \frac{1}{2} [(\underline{x}_1 - \underline{\mu}_1)' \Omega_{11}(\underline{x}_1 - \underline{\mu}_1) + 2(\underline{x}_1 - \underline{\mu}_1)' \Omega_{12}(\underline{x}_2 - \underline{\mu}_2)]$$

$$= c - \frac{1}{2} (\underline{x}_1 - \underline{\mu}_1^*)' \Omega_{11} (\underline{x}_1 - \underline{\mu}_1^*), \text{ where } \underline{\mu}_1^* \text{ satisfies}$$

$$2\underline{x}_1' \Omega_{11} \underline{\mu}_1^* = 2\underline{x}_1' \Omega_{11} \underline{\mu}_1 - 2\underline{x}_1' \Omega_{12}(\underline{x}_2 - \underline{\mu}_2)$$

$$= 2\underline{x}_1' \Omega_{11} (\underline{\mu}_1 - \Omega_{11}^{-1} \Omega_{12}(\underline{x}_2 - \underline{\mu}_2))$$

so $\underline{x}_1 | \underline{x}_2$ follows a multivariate Gaussian distribution with

mean $\underline{\mu}_1 - \Omega_{11}^{-1} \Omega_{12}(\underline{x}_2 - \underline{\mu}_2)$ and var Ω_{11} . It is also the distribution of a regression target of model

$$y = \left(\underbrace{\underline{\mu}_1}_{\text{Intercept}} + \underbrace{\Omega_{11}^{-1} \Omega_{12}}_{\text{slope}} \underline{x}_2 \right) + \epsilon, \quad \epsilon \sim N(0, \Omega_{11})$$

Intercept

slope

where randomness comes from

Multiple regression: + three classical principles for inference.

(A) Probably the coolest aspect of linear regression!

Least squares: $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \underline{x}_i' \beta)^2 \right\}$

Maximum of Gaussianity:

$$\begin{aligned} \hat{\beta} &= \arg \max_{\beta \in \mathbb{R}^p} \left\{ \prod_{i=1}^n p(y_i | \beta, \sigma^2) \right\} = \arg \max_{\beta \in \mathbb{R}^p} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \underline{x}_i' \beta)^2 \right\} \\ &= \arg \min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \underline{x}_i' \beta)^2 \right\}, \text{ which reduces to Least Squares problem.} \end{aligned}$$

Method of moments

Denote matrix $J_n = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n'$, then $(I - J_n)$ is the "centralize operator" i.e., $(I - J_n)\underline{x} = (x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x})'$. It's not hard to see that matrix $(I - J_n)$ is symmetric and idempotent, i.e.,

$$\begin{cases} (I - J_n)' = (I - J_n) \\ (I - J_n)^m = (I - J_n) \end{cases}$$

Denote error by vector \underline{e} , $\underline{e} = \{e_i\}_{n \times 1}$, $e_i = y_i - \underline{x}_i' \beta$,

or $\underline{e} = \underline{y} - \underline{X}\beta$, \underline{x}_i' is the i -th column of \underline{X} . Then,

$\text{cov}(\underline{x}, \underline{e}) = 0 \Leftrightarrow \underline{X}'(I - J_n)'(I - J_n)\underline{e} = 0$, by the symmetry and idempotency of $(I - J_n)$, we have

$$\text{cov}(\underline{x}, \underline{e}) = 0 \Leftrightarrow \underline{X}'(I - J_n)'\underline{e} = 0$$

$$\Leftrightarrow \begin{bmatrix} x_{11} - \bar{x}_1 & x_{21} - \bar{x}_1 & \dots & x_{n1} - \bar{x}_1 \\ x_{12} - \bar{x}_2 & x_{22} - \bar{x}_2 & \dots & x_{n2} - \bar{x}_2 \\ \vdots & \vdots & \ddots & \vdots \\ x_{1p} - \bar{x}_p & x_{2p} - \bar{x}_p & \dots & x_{np} - \bar{x}_p \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (*)$$

if we use standardized data, i.e., $\bar{x}_1 = \bar{x}_2 = \dots = \bar{x}_p = \bar{y} = 0$, then

Eq. (*) is equivalent to Least Squares statement.

$$(B) \arg \min_{\beta \in \mathbb{R}^p} \left\{ \sum_{i=1}^n w_i (y_i - x_i' \beta)^2 \right\}$$

$$= \arg \max_{\beta \in \mathbb{R}^p} \exp \left\{ - (y - X\beta)' W (y - X\beta) \right\}, \quad W = \text{diag}\{w_1, \dots, w_p\}$$

The RHS can be interpreted as maximizing likelihood of a heteroscedastic Gaussian with precision matrix W , and

$$W = \Sigma^{-1} = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_p^2)^{-1} = \text{diag}(1/\sigma_1^2, \dots, 1/\sigma_p^2)$$

Quantifying uncertainty: some basic frequentist ideas.

In linear regression Notation: Both y and β are vectors, X is data matrix.

$$(A) \hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} L \quad \text{where } L = (y - X\beta)'(y - X\beta), \text{ then}$$

$$\nabla_{\beta} L|_{\beta=\hat{\beta}} = -2X'y + 2X'X\hat{\beta} = 0 \Rightarrow \hat{\beta} = (X'X)^{-1}X'y$$

Since $y = X\beta + \epsilon \sim N(X\beta, \sigma^2 I)$, then

$$\hat{\beta} = (X'X)^{-1}X'y \sim N((X'X)^{-1}X'X)\beta, (X'X)^{-1}X'\sigma^2 I X(X'X)^{-1}) = N(\beta, (X'X)^{-1}\sigma^2)$$

(B) Apparently, the standard errors for each β_j is just the square root of the (j, j) element in $(X'X)^{-1}\sigma^2$. By the relationship between adjoint and inverse of a matrix, if mat A is invertible then

$$A^{-1} = \frac{\text{Adj. } A}{|A|}$$

so the (j, j) element in A^{-1} is $\text{Adj. } A [j, j] / |A|$.

Going back to our problem,

$$\text{var}(\beta_j) = \sigma^2 \frac{\text{Adj. } (X'X) [j, j]}{|X'X|} = \sigma^2 \frac{|M_{jj}| \cdot (-1)^{(j+j)}}{|X'X|} = \frac{\sigma^2 |M_{jj}|}{|X'X|}, \text{ when}$$

M_{jj} is the minor of $X'X$ that removes the j -th row and j -th column in $(X'X)$. Denote the i -th column in X by x_i , then $(X'X) [i, j] = x_i'x_j$,

therefore by removing the j -th row and column from $(X'X)$, we are virtually doing $(X_{-j}'X_{-j})$, where X_{-j} denotes the remaining cols in X after removing column j

$$\text{As a result, } \text{var}(\beta_j) = \sigma^2 \cdot \frac{|X_{-j}'X_{-j}|}{|X'X|}$$

Propagating uncertainty.

$$(A) \quad \underline{\theta} = (\theta_1, \dots, \theta_p)', \quad f(\underline{\theta}) = \theta_1 + \theta_2 \Leftrightarrow f(\underline{\theta}) = [1, 1, 0, \dots, 0] \underline{\theta},$$
$$\text{then } \text{var}(f(\hat{\underline{\theta}})) = [1, 1, 0, \dots, 0] \hat{\Sigma} \begin{bmatrix} 1 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \hat{\Sigma}_{11} + 2\hat{\Sigma}_{12} + \hat{\Sigma}_{22}$$

Similarly, if $f(\underline{\theta}) = [1, 1, \dots, 1] \underline{\theta}$, then

$$\text{var}(f(\hat{\underline{\theta}})) = [1, 1, \dots, 1] \hat{\Sigma} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \sum_{i=1}^p \sum_{j=1}^p \hat{\Sigma}_{ij} \quad (i, j \text{ are allowed to be identical})$$

(B) What (A) tells is how to calculate var of $f(\underline{\theta})$ when f itself is a linear function w.r.t. $\underline{\theta}$. So the idea of calculating var of $f(\underline{\theta})$ when f is not linear is to make a linear approximation of f and calculate the var of this linear approximation.

$$f(\hat{\underline{\theta}}) \approx f(\underline{\theta}) + (\nabla_{\underline{\theta}} f)'(\hat{\underline{\theta}} - \underline{\theta}) \quad (1\text{-st order Taylor approximation})$$

$$\begin{aligned} \text{Then } \text{var}(f(\hat{\underline{\theta}})) &\approx \text{var}[f(\underline{\theta}) + (\nabla_{\underline{\theta}} f)'(\hat{\underline{\theta}} - \underline{\theta})] \\ &= (\nabla_{\underline{\theta}} f)' \text{var}(\hat{\underline{\theta}} - \underline{\theta}) (\nabla_{\underline{\theta}} f) \end{aligned}$$

Caveats:

- (i) To make linear approximation make sense, f should not change violently (have a large second order derivative) w.r.t. $\underline{\theta}$; and $\hat{\underline{\theta}}$ should not be far from $\underline{\theta}$, so $\hat{\underline{\theta}}$ almost must be an unbiased estimator for $\underline{\theta}$.