

⑤

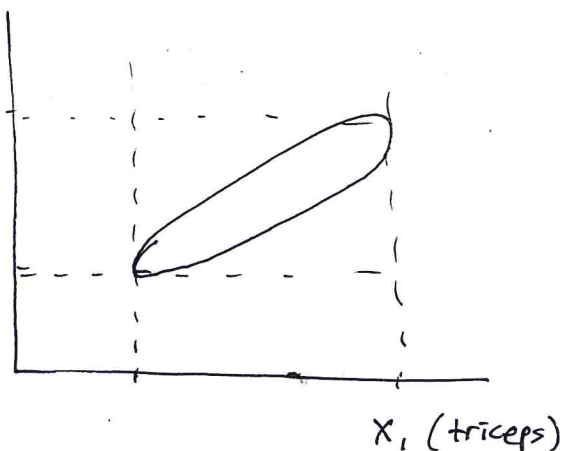
Multicollinearity exists when input variables are highly correlated among themselves

example: $X_1 = \text{triceps}$, $X_2 = \text{thigh}$, $X_3 = \text{midarm}$, $y = \text{body fat}$
(X_1, X_2, X_3 are easy measurements, y is cumbersome and expensive)

correlation matrix $r_{12} = .924$, $r_{13} = .458$, $r_{23} = .085$

$r_{y1} = .843$, $r_{y2} = .878$, $r_{y3} = .142$

X_2 (thigh)



[Input space is multi-dimensional,
not necessarily rectangular,

[Analysis must include an
understanding of how input variables
are related

Thoughts:

(1) Interpretation of a regression coefficient depends on which other inputs are in the model

examples → $y = \text{body fat}$, $X_1 = \text{weight}$, $X_2 = \text{abdomen size}$

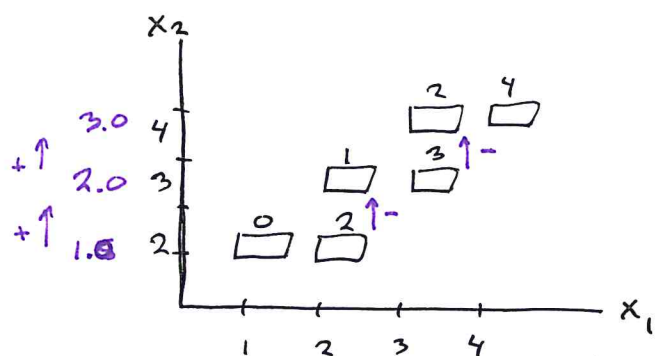
→ $y = \text{exam score}$, $X_1 = \text{difficulty}$, $X_2 = \text{ability}$

→ $y = \text{school quality}$, $X_1 = \text{money spent}$, $X_2 = \text{teacher quality}$

→ $y = \text{wins}$, $X_1 = \text{shots on goal}$, $X_2 = \text{goals scored}$

(6)

simple example : $E(Y) = 2X_1 - X_2$, $r_{12} > 0$.



X_2 has a negative partial effect,
given a fixed level of X_1

X_2 has a positive marginal effect,
not accounting for the level of X_1

example : $\hat{y} = 117.085 + 4.334 X_1 - 2.857 X_2 - 2.186 X_3$
(recall that $r_{Y2} > 0$)

(2) If an input is highly correlated with response,
additional inputs are limited in how much new information
can be provided.

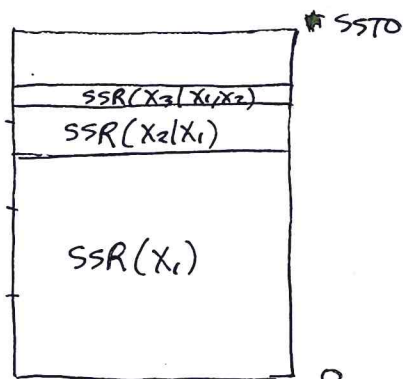
example : $SSR(X_1) = 352.27$ $R^2(X_1) = .711$
 $SSR(X_2|X_1) = 33.17$ $R^2(X_1, X_2) = .778$
 $SSR(X_3|X_1, X_2) = 11.55$ $R^2(X_1, X_2, X_3) = .801$
 (recall that $r_{Y1} = .843$
 $r_{Y2} = .878$)

$$F^* = 21.52 (p = .000)$$

$$t_1^* = 1.437 (p = .170)$$

$$t_2^* = -1.106 (p = .285)$$

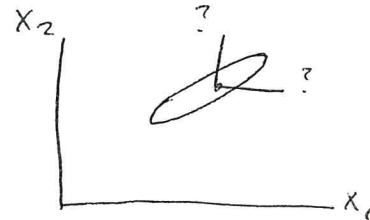
$$t_3^* = -1.370 (p = .190)$$



⑦

(3) Estimated regression coefficients have large sampling variance.

idea: We have little information on the effect of X_1 (X_2) with X_2 (X_1) held fixed.



example:

Model	b_1 (SE(b_1))	b_2 (SE(b_2))
X_1	.8572 (.1288)	—
X_2	—	.8565 (.1100)
X_1, X_2	.2224 (.3034)	.6594 (.2912)
X_1, X_2, X_3	4.334 (3.016)	-2.857 (2.582)

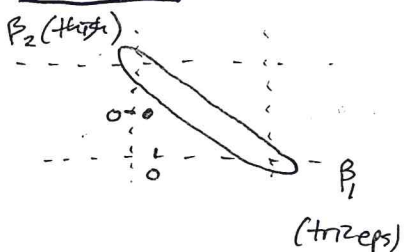
(4) Prediction and mean estimation (fitted regression function) is not affected by multicollinearity within the input space

example:

X_1	X_2	X_3	CI for μ_h
25	.	.	[18.6, 21.3]
25	50	.	[18.0, 20.7]
25	50	29	[17.9, 20.5]

(5) Regression coefficient estimates are highly correlated when input variables are highly correlated.

example:



$\beta_1 > b_1$ is compatible with $\beta_2 < b_2$

$\beta_1 < b_1$ is compatible with $\beta_2 > b_2$.

$(0,0)$ is contained in $CI_{\beta_1} \times CI_{\beta_2}$

but $(0,0)$ not contained in $CE(\beta_1, \beta_2)$

illustrates the need to think in higher dimensions

8

Investigating relationships in higher dimensions requires higher level statistical methods, such as regression analysis. Two-dimensional methods and graphs are insufficient.

example: $(P_1) X_1 = \text{triceps}, (P_2) X_2 = \text{thigh}, (W) X_3 = \text{midarm}$

$$r_{1W} = .458, r_{2W} = .085$$

(moderate, small correlations between X_1 and X_2)

$$r_{1W}^2 = .2096, t_1^* = 2.185, p = .04$$

$$(F_1^* = 4.77)$$

$$r_{2W}^2 = .0072, t_2^* = 0.361, p = .72$$

$$(F_2^* = 0.13)$$

effect marginal tests
 $M_2: W = \beta_0 + \beta_2 X_2 + \epsilon$
 $M_0: W = \beta_0 + \epsilon$

The observed data is compatible with the reduced model. It is not necessary to add thigh measurement to the no effects model for predicting midarm measurement.

Now look at what happens when both X_1, X_2 are used to model W .

$$r_{3,12}^2 = .99, F_{12}^* = 880.7, p_{12} = .000$$

(very high correlation between X_3 and (X_1, X_2))

$$t_{1|2}^* = 41.82, t_{2|1}^* = -37.26$$

$$(p = .000)$$

$$(p = .000)$$

$$(F_{1|2}^* = 1748.6)$$

$$(F_{2|1}^* = 1388.6)$$

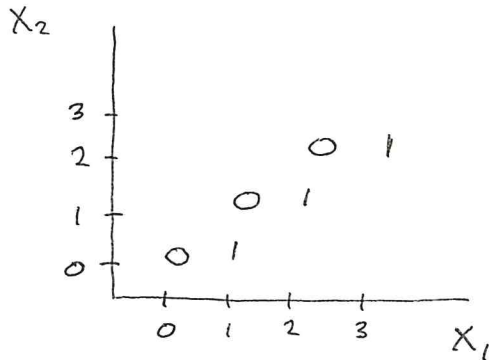
partial effect test
 $M_{12}: W = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$
 $M_1: W = \beta_0 + \beta_1 X_1 + \epsilon$

The observed data is not compatible with the reduced model.

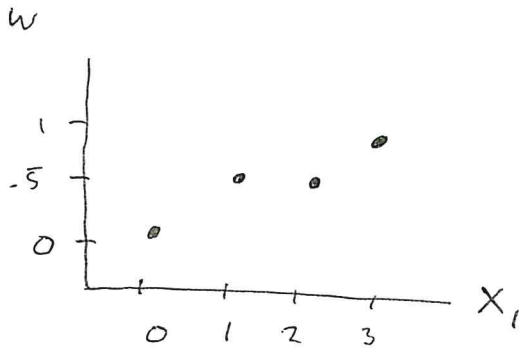
We accept the addition of thigh measurement as a predictor for midarm to the model which already includes triceps.

⑧ (b)

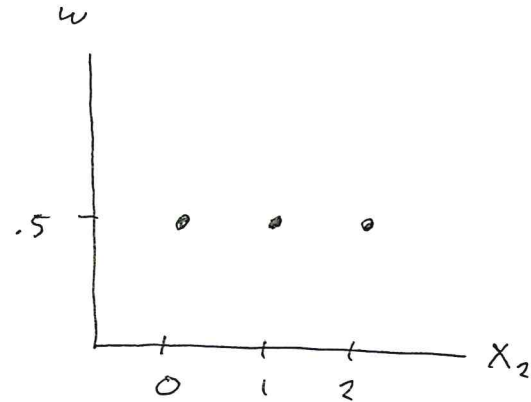
Simple example: $w = X_1 - X_2$, X_1, X_2 are positively correlated



$\text{Corr}(w, (X_1, X_2))$ is large



$\text{Corr}(w, X_1)$ is moderate



$\text{Corr}(w, X_2)$ is small

example: $\hat{\text{midarm}} = 62.33 + 1.88 \text{ triceps} - 1.6 \text{ thigh}$

①

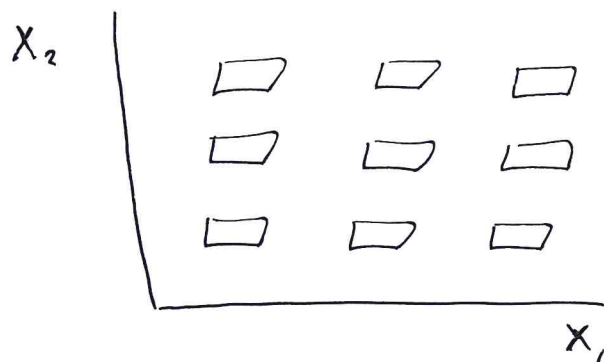
Orthogonal Design Matrix :

Def:

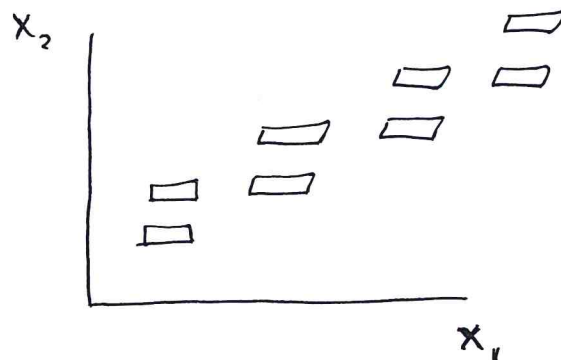
[A design is orthogonal if $X'X$ is diagonal.]

idea:

Orthogonal Design



Correlated Inputs



marginal effect of X_1 : average \updownarrow , take differences

partial effect of X_1 : differences \leftrightarrow , then average

partial effect of X_1 is determined for constant X_2 conditions
under orthogonality, ^{levels} ~~of X_1~~ are the same across X_2 .

example : study of coal conservation into oils, Canadian lignite

y = conversion percentage

X_1 = temperature, X_2 = molar ratio ($\frac{CO}{H_2}$)
 (380, 460) °C (1/4, 3/4)

X_3 = pressure, X_4 = contact time
 (7.1, 11.1) MPa (10, 50) min

$n = 16$, $r_{X_i, X_k} = 0$ (2^4 Design)

②

To get $X'X$ diagonal, we use coded variables:

$$X_i = \begin{cases} +1, & \text{input at high level} \\ -1, & \text{input at low level} \end{cases}$$

Note that $\text{Cov}(\underline{b}) = \sigma^2 (X'X)^{-1}$ is diagonal. $\left(\text{Cov}(b_\ell, b_k) = 0 \text{ for all } \ell \neq k \right)$

Also, $\text{SSR}(X_\ell | X_k, k \neq \ell) = \text{SSR}(X_\ell)$

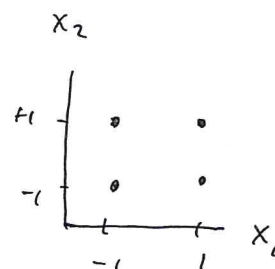
(Thus, there is no ambiguity in defining the effect of an input)

For orthogonal designs, regression coefficient estimates and variation explained by an input do not depend on which other inputs are included in the model.

example: 2^k design, $b_\ell = \frac{1}{2}(\bar{Y}_{h\ell} - \bar{Y}_{l\ell})$

$$t_\ell^* = \frac{b_\ell}{\text{SE}(b_\ell)}$$

$$\text{SE}(b_\ell) = \sqrt{\frac{\text{MSE}}{n}}$$



$b_1 = \frac{5.0}{1}, b_2 = \frac{3.15}{1}, b_3 = \frac{3.975}{1}, b_4 = \frac{1.862}{1}, \text{SE}(b_\ell) = \frac{1.408}{1}$

$(p_1 = .005), (p_2 = .047), (p_3 = .017), (p_4 = .213)$

selected model includes X_1, X_2, X_3 .

$b_1 = \frac{5.0}{1}, b_2 = \frac{3.15}{1}, b_3 = \frac{3.975}{1}, \text{SE}(b_\ell) = \frac{1.451}{1}$ (MSE increases when X_4 is removed)

$(p_2 = .051)$

Do not get stuck in dichotomous thinking in model selection.