

Modern Data Science Methods for Educational Research

R for Data Analysis in Educational Research

Data Analysis

ผศ.ดร.สิวะโชติ ศรีสุทธียากร
อ.ดร.ประภาศิริ รัชประภาพรกุล

ภาควิชาวิจัยและจิตวิทยาการศึกษา
คณะครุศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย

February 12, 2023



KRUROO
EDU @ CHULALONGKORN

Modern Data Science Methods for Educational Research

สำรวจการแจกแจงของตัวแปร : สถิติพื้นฐาน

```
1 summary(dat)
```

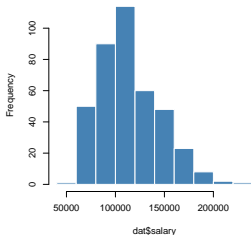
rank	discipline	yrs.since.phd	yrs.service
AssocProf: 64	A:181	Min. : 1.00	Min. : 0.00
AsstProf : 67	B:216	1st Qu.:12.00	1st Qu.: 7.00
Prof :266		Median :21.00	Median :16.00
		Mean :22.31	Mean :17.61
		3rd Qu.:32.00	3rd Qu.:27.00
		Max. :56.00	Max. :60.00
salary			
Min. : 57800			
1st Qu.: 91000			
Median :107300			
Mean :113706			
3rd Qu.:134185			

สำรวจการแจกแจงของตัวแปร : ตัวแปรต่อเนื่อง

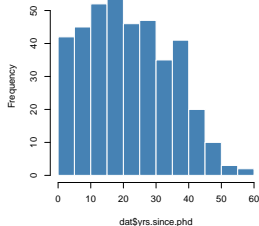
```
1 par(mfrow=c(2,3))
2 hist(dat$salary, col="steelblue", border="white")
3 hist(dat$yrs.since.phd, col="steelblue", border="white")
4 hist(dat$yrs.service, col="steelblue", border="white")
5 boxplot(dat$salary, horizontal = TRUE)
6 boxplot(dat$yrs.since.phd, horizontal = TRUE)
7 boxplot(dat$yrs.service, horizontal = TRUE)
```

สำรวจการแจกแจงของตัวแปร : ตัวแปรต่อเนื่อง

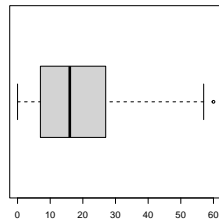
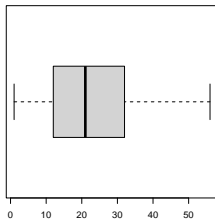
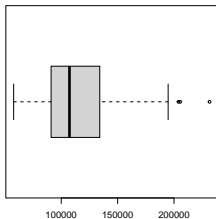
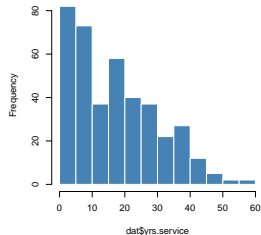
Histogram of dat\$salary



Histogram of dat\$yrs.since.phd



Histogram of dat\$yrs.service



สำรวจการแจกแจงของตัวแปร : ตัวแปรจัดประเภท

```
1 tab_rank <- table(dat$rank)
2 tab_rank
```

AssocProf	AsstProf	Prof
64	67	266

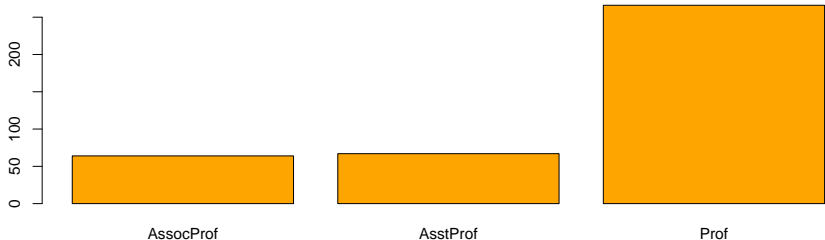
```
1 tab_discipline <- table(dat$discipline)
2 tab_discipline
```

A	B
181	216

```
1 tab_sex <- table(dat$sex)
2 tab_sex
```

สำรวจการแจกแจงของตัวแปร : ตัวแปรจัดประเภท

```
1 barplot(tab_rank, col = "orange")
```



วัตถุประสงค์ลักษณะที่ 1

เพื่อเปรียบเทียบค่าเฉลี่ยเงินเดือนของอาจารย์มหาวิทยาลัย ระหว่าง
กลุ่มอาจารย์ที่มี ตำแหน่งทางวิชาการ สาขาวิชา และเพศแตกต่างกัน

- ▶ **สำรวจความสัมพันธ์ (วิเคราะห์เปรียบเทียบ) เงินเดือน**
จำแนกตามตำแหน่งวิชาการ สาขาวิชา และเพศ
- ▶ **สำรวจความสัมพันธ์ระหว่างตำแหน่งวิชาการ สาขาวิชา และเพศ**

Modern Data Science Methods for Educational Research

	discipline	mean	sd	min	max
	<fct>	<dbl>	<dbl>	<int>	<int>
1	Pure Science	108548.	30538.	57800	205500
2	Applied Science	118029.	29459.	67559	231545

สรุปผลการสำรวจ

- ▶ เงินเดือนของอาจารย์มหาวิทยาลัย มีแนวโน้มแตกต่างกันได้ตามตำแหน่งวิชาการ สาขาวิชา และเพศ
- ▶ มีแนวโน้มว่า จะอิทธิพลปฏิสัมพันธ์ของตำแหน่งวิชาการ และ เพศ ต่อเงินเดือนอาจารย์

Note: ตารางแจกแจงความถี่สองทาง

```
1 table(dat$rank, dat$discipline)
```

	Pure Science	Applied Science
AsstProf	24	43
AssocProf	26	38
Prof	131	135

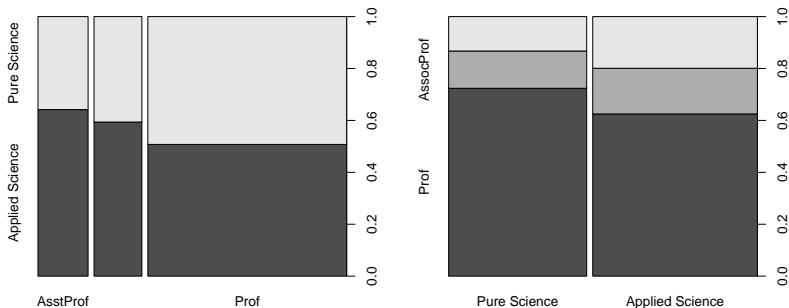
```
1 table(dat$rank, dat$sex)
```

	Female	Male
AsstProf	11	56
AssocProf	10	54
Prof	18	248

```
1 table(dat$discipline, dat$sex)
```


Note: Mosaic plot

```
1 par(mfrow=c(1,2))  
2 plot(dat$rank, dat$discipline, xlab = " ", ylab = " ")  
3 plot(dat$discipline, dat$rank, xlab = " ", ylab = " ")
```



Modelling สำหรับวัตถุประสงค์ลักษณะแรก

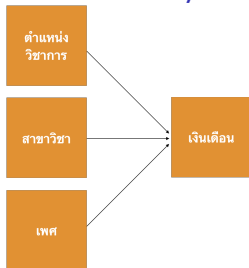
Job 1: independent-sample t-test



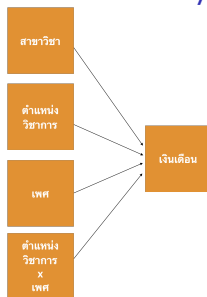
Job 2: One-Way ANOVA



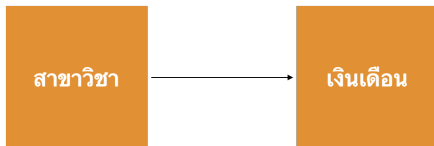
Job 3: three-ways ANOVA



Job 4: Three-Way ANOVA 2



Job1: Independent-sample t-test



สมมุติฐานการทดสอบ

$$H_0 : \mu_{applied} \leq \mu_{pure}$$

$$H_1 : \mu_{applied} > \mu_{pure}$$

Job1: Independent-sample t-test

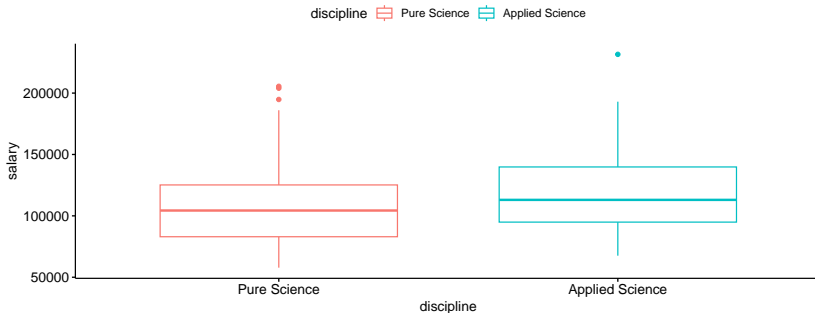
```
1 dat %>% group_by(discipline) %>%  
2   summarise(mean = mean(salary),  
3             sd = sd(salary),  
4             min = min(salary),  
5             max = max(salary))
```

```
# A tibble: 2 x 5
```

	discipline	mean	sd	min	max
	<fct>	<dbl>	<dbl>	<int>	<int>
1	Pure Science	108548.	30538.	57800	205500
2	Applied Science	118029.	29459.	67559	231545

Job1: Boxplot using ggpubr

```
1 #install.packages("ggpubr")  
2 library(ggpubr)  
3 ggboxplot(data = dat, x = "discipline", y = "salary",  
4           color = "discipline")
```



Job1: Independent-sample t-test

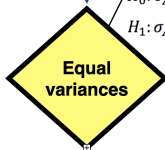
ขั้นตอนการวิเคราะห์

Independent sample t-test

ตรวจสอบได้ด้วย Levene's test

$$H_0: \sigma_{ACH_{PBL}}^2 = \sigma_{ACH_{Lec}}^2$$

$$H_1: \sigma_{ACH_{PBL}}^2 \neq \sigma_{ACH_{Lec}}^2$$



Equal
variances

Yes.
(Don't reject H_0
of Levene's test.)

Pooled variances t-

$$t^* = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}; df = n_1 + n_2 - 2$$

โดยที่ S_p คือ pooled standard deviation

$$S_p = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}}$$

No. (Reject H_0 of Levene's test)

Separated variances t-test

$$t^* = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

$$df = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\left(\frac{S_1^2}{n_1}\right)^2 \left(\frac{1}{n_1 - 1}\right) + \left(\frac{S_2^2}{n_2}\right)^2 \left(\frac{1}{n_2 - 1}\right)}$$

Job1: Levene's Test for equality of variances

```

1 #install.packages("car")
2 library(car)
3 leveneTest(salary ~ discipline, data = dat)

```

Levene's Test for Homogeneity of Variance (center = median)

	Df	F value	Pr(>F)
group	1	0.0458	0.8306
	395		

Job1: Syntax for t-test

```
1 # two-sided test
2 t.test(salary ~ discipline ,data=dat,
3       var.equal = TRUE)
4 # one-sided test
5 t.test(salary ~ discipline, data=dat,
6       var.equal = FALSE,
7       alternative = "greater")
```


Job1: Independent-sample t-test output

```
1 t.test(salary ~ discipline, data=dat,
2       var.equal = FALSE,
3       alternative = "less")
```

Welch Two Sample t-test

data: salary by discipline

t = -3.1306, df = 377.83, p-value = 0.00094

alternative hypothesis: true difference in means between gr

95 percent confidence interval:

-Inf -4487.034

sample estimates:

mean in group Pure Science mean in group Applied Science

108548.4

118028.7

Job1: Assumptions check

independent sample t-test มีข้อตกลงเบื้องต้นที่สำคัญดังนี้

- ▶ Independence
- ▶ Normality
- ▶ Homogeneity of variances

Job1: Assumptions check

```
1 shapiro.test(dat$salary)
```

Shapiro-Wilk normality test

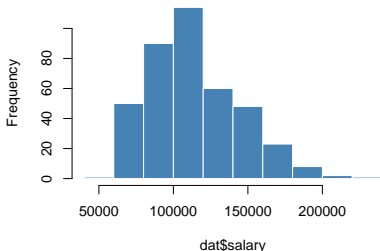
data: dat\$salary

W = 0.95988, p-value = 6.076e-09

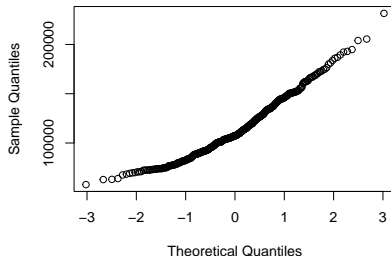
Job1: Assumptions check

```
1 par(mfrow=c(1,2))
2 hist(dat$salary, col="steelblue", border = "white")
3 qqnorm(dat$salary)
```

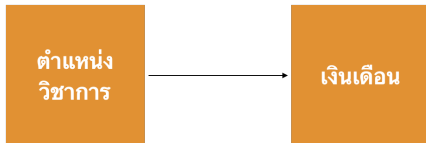
Histogram of dat\$salary



Normal Q-Q Plot



Job2: Independent-sample F-test (One-Way ANOVA)



```
1 dat %>% group_by(rank) %>%  
2   summarise(mean = mean(salary),  
3             sd = sd(salary),  
4             min = min(salary),  
5             max = max(salary))%>%  
6   arrange(desc(mean))
```

Job2: Independent-sample F-test (One-Way ANOVA)

```
# A tibble: 3 x 5
```

	rank	mean	sd	min	max
	<fct>	<dbl>	<dbl>	<int>	<int>
1	Prof	126772.	27719.	57800	231545
2	AssocProf	93876.	13832.	62884	126431
3	AsstProf	80776.	8174.	63100	97032

Job2: Calculate F-test for Overall hypothesis

$$H_0 : \mu_1 = \mu_2 = \mu_3$$

$$H_1 : \text{not } H_0$$

```
1 # calculate analysis of variance
2 fit <- aov(salary ~ rank, data= dat)
3 summary(fit)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
rank	2	1.432e+11	7.162e+10	128.2	<2e-16 ***
Residuals	394	2.201e+11	5.586e+08		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

ผลการวิเคราะห์ข้างต้นสรุปว่า ...

Job2: Multiple comparison

การทดสอบเพื่อเปรียบเทียบค่าเฉลี่ยรายคู่มีหลายวิธีการ

- ▶ TukeyHSD()
- ▶ `pairwise.t.test()`
- ▶ ScheffeTest()

Job3: Summary Stat

```
# A tibble: 7 x 4
# Groups:   ind [3]
  ind      value      mean    sd
<chr>    <chr>    <dbl>  <dbl>
1 discipline Applied Science 118029. 29459.
2 discipline Pure Science   108548. 30538.
3 rank      AssocProf      93876. 13832.
4 rank      AsstProf       80776.  8174.
5 rank      Prof          126772. 27719.
6 sex       Female         101002. 25952.
7 sex       Male           115090. 30437.
```


Job3: Calculate Multiple Comparison

```
1 TukeyHSD(fit3)
```

Tukey multiple comparisons of means
95% family-wise confidence level

```
Fit: aov(formula = salary ~ rank + discipline + sex, data =
```

\$rank

	diff	lwr	upr	p adj
AssocProf-AsstProf	13100.45	3790.209	22410.70	0.0029189
Prof-AsstProf	45996.12	38715.029	53277.22	0.0000000
Prof-AssocProf	32895.67	25479.514	40311.83	0.0000000

\$discipline

diff lwr upr p a

Job4: Summary Stat

```
1 dat %>% group_by(rank, sex) %>%
2   summarise(mean = mean(salary),
3             sd = sd(salary))
```

```
# A tibble: 6 x 4
```

```
# Groups:    rank [3]
```

	rank	sex	mean	sd
	<fct>	<fct>	<dbl>	<dbl>
1	AsstProf	Female	78050.	9372.
2	AsstProf	Male	81311.	7901.
3	AssocProf	Female	88513.	17965.
4	AssocProf	Male	94870.	12891.
5	Prof	Female	121968.	19620.
6	Prof	Male	127121.	28214.

