

Moderation & Mediation

Rui Zhou
Spring 2023

目录

1 这章在讲什么	1
2 Moderation	2
2.1 Fit Moderation Model	3
2.2 After Significant	5
2.3 Appendix for Moderation	6
3 Mediation	8
3.1 Fit Mediation Model	9
3.2 Mediation Effect	11
3.3 Appendix for Mediation	11
4 容易踩的坑	14

1 这章在讲什么

调节 (moderation) 和中介 (mediation) 解决的是两个**不同**问题，初学者经常混。一句话区分：

- **调节** $X \xrightarrow{Z} Y$: X 对 Y 的影响**取决于第三个变量 Z** (年纪小的时候社会支持作用大，年纪大就没那么明显)。回归形式: $Y = a + b_1X + b_2Z + b_3(X \cdot Z) + \varepsilon$, 关键看 b_3 (交互项系数) 的显著性。

- **中介** $X \rightarrow M \rightarrow Y$: X 通过一个**中间变量** M 影响 Y (学习时间 \rightarrow 自我效能感 \rightarrow 学业成绩)。三步走: $Y \sim X$ (总效应)、 $M \sim X$ (a 路径)、 $Y \sim X + M$ (b 路径 + 直接效应)。

两个核心数学量

调节的 simple slope: 当 Z 取某个具体值 z_0 时, X 对 Y 的斜率为 $b_1 + b_3 z_0$ 。常规做法是在 Z 的均值、 $\bar{Z} \pm 1 \text{sd}$ 三个点上画斜率, 看哪些斜率显著、哪些不显著。

中介的间接效应 $a \cdot b$: X 经过 M 影响 Y 的部分。检验显著性首选**自助法置信区间** (bootstrap) 而不是 Sobel 检验——因为 $a \cdot b$ 的抽样分布并不正态, Sobel 功效较低、保守。Hayes 的 PROCESS 宏 / R 的 `mediation` 包默认就是 bootstrap。

中心化变量 (centering) 在调节分析里几乎是必做的: 把 X 和 Z 各自减去均值再相乘构造交互项, 可以缓解 $X \cdot Z$ 与 X 、 Z 本身的多重共线性, 让主效应 b_1 、 b_2 解读为「在另一个变量取均值时的效应」。

2 Moderation

In statistics, moderation refers to a type of **interaction effect** between two or more variables in a statistical model. Specifically, it occurs when the strength or direction of the relationship between two variables (the predictor and the outcome variable) changes depending on the level of a third variable (the moderator). Moderator variable is a variable whose different values determine the nature of the relationship between two other variables.

To test for moderation in a statistical model, researchers typically include an **interaction term** between the *predictor* and the *moderator* variable. If the interaction term is statistically significant, it suggests that the relationship between the predictor and outcome variable is different at different levels of the moderator variable. Moderation in model is an extension for two-way ANOVA, without limits to the scale of data (categorical or continuous).

Moderation analysis is important in statistics because it allows researchers

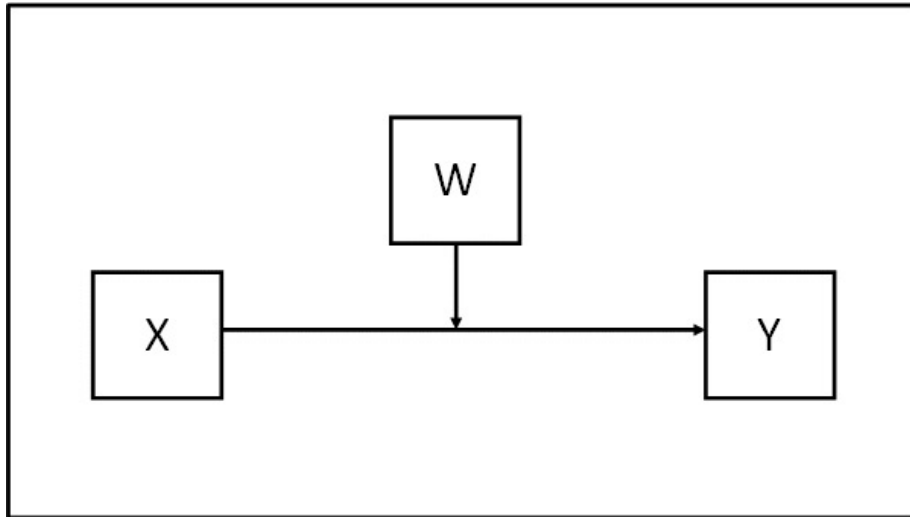


图 1: Moderation

to examine the conditions under which a relationship between two variables is stronger or weaker. This information can be useful for developing interventions or identifying subgroups of individuals who may benefit from different types of treatments or interventions.

2.1 Fit Moderation Model

2.1.1 NULL Model

The data in “moderation_data.sav” contains information about loyalty to marriage, relationship between partners, and their age. We’re interested in the relationship between age and loyalty, and how the relationship in marriage moderates it.

```
# moderation data  
currentData <- bruceR::import("moderation_data.sav")  
attach(currentData)
```

Note that when an interaction term is introduced into the model, better

to centralize the variable for convenience of interpretation! (Even without centered independent variables, the significance output won't change.)

```
C_Relationship = Relationship - mean(Relationship)
C_age_con = age_con - mean(age_con)
```

The NULL model here is without the interaction term, but to inspect age and Relationship through MLR.

```
model_0 <- lm(Loyalty ~ C_Relationship + C_age_con)
bruceR::GLM_summary(model_0)
```

2.1.2 Moderation Model

Then, add the interaction of moderator and predictor into the MLR model. (Moderation in essence is a MLR with interaction terms.)

```
model_moderation <- lm(Loyalty ~ C_Relationship * C_age_con)
bruceR::GLM_summary(model_moderation)
```

After introducing the interaction term into our model, four things are of great interest.

- Significance output of interaction term.
- Significance output of original terms in null model.
- Adjusted R^2 compared to null model.
- VIF, in case of multicollinearity.

The improvement of R^2 indicates that, even after considering the added complexity of model, the explanatory power is greater than the null.

2.1.3 Model Comparison

The most natural impulse is to compare the moderation model against the null one. The difference lies in the interaction term. In fact, the output comes down right to the interaction term.

```
anova(model_0, model_moderation)
```

2.2 After Significant

After seeing the moderation (interaction) is significant, we need to check how it exerts its effect by **simple main effect**, because a significant moderation term will affect your interpretation of main effect.

We're to compute marginal means, where the data is divided into 3×3 categories using $Mean \pm SD$, i.e., with 3 X levels and 3 M levels.

```
library(emmeans)
m_Relationship <- mean(Relationship, na.rm = TRUE)
sd_Relationship <- sd(Relationship, na.rm = TRUE)
m_Age <- mean(age_con, na.rm = TRUE)
sd_Age <- sd(age_con, na.rm = TRUE)

# use un-centered data to see how the raw looks like in slope analysis.
model_moderation <- lm(Loyalty ~ Relationship * age_con)

# compute Estimated Marginal Means
emm <- emmeans(model_moderation, ~Relationship * age_con, cov.keep = 3, at = list(age_con = c(
  sd_Age, m_Age, m_Age + sd_Age), Relationship = c(m_Relationship - sd_Relationship,
  m_Relationship, m_Relationship + sd_Relationship)), level = 0.95)
summary(emm)
```

We then do the slope analysis, similar to simple main effect.

```
simpleSlope <- emtrends(model_moderation, pairwise ~ age_con, var = "Relationship",
  cov.keep = 3, at = list(age_con = c(m_Age - sd_Age, m_Age, m_Age + sd_Age)),
  level = 0.95)
summary(simpleSlope)
```

Slopes can be visualized.

```
emmip(model_moderation, age_con ~ Relationship, cov.keep = 3, at = list(Relationship =
  sd_Relationship, m_Relationship, m_Relationship + sd_Relationship), age_con = c(m_A
  sd_Age, m_Age, m_Age + sd_Age)), CIs = TRUE, level = 0.95, position = "jitter")
```

2.3 Appendix for Moderation

2.3.1 Uncentered Model

One plausible reason for centralization is interpretation, and the other is possible collinearity.

```
model_nonc = lm(Loyalty ~ Relationship + age_con + Relationship * age_con)
bruceR::GLM_summary(model_nonc)
```

Even though the uncentered model reports the same results with respect to significance, the VIF bombed!

2.3.2 bruceR::GLM_summary()

By the way, why do I use `bruceR::GLM_summary(model)` here? Because it will return well-rounded report with standardized and non-standardized coefficients, 95% CI for coefficients, significance output for variables, necessary report of the model. Otherwise, if I'm to get the 95% CI of each slope coefficients, I have to construct a new command like:

```
cbind(coef(model_moderation), confint(model_moderation, level = 0.95))
```

2.3.3 Residual Plots

```
res.std <- rstandard(model_moderation)
plot(res.std, ylab = "Standardized Residuals")
car::residualPlots(model_moderation)

library(ggplot2)
ggplot(as.data.frame(res.std), aes(sample = res.std)) + geom_qq() + geom_qq_line()
```

Lastly,

```
detach(currentData)
```

2.3.4 More on Moderation

There are several limitations to moderation analysis that researchers should be aware of:

1. Causality cannot be inferred: Moderation analysis can only identify statistical relationships between variables, and cannot establish causality. To establish causality, researchers need to use experimental designs or quasi-experimental designs that include random assignment.
2. Measurement error: Moderation analysis assumes that all variables are measured without error. However, in practice, measurement error is often present, which can lead to biased estimates of the moderation effect.
3. Limited generalizability: Moderation effects may be specific to the sample, context, or time period in which the study was conducted, and may not generalize to other populations or contexts.
4. Sample size requirements: Moderation analysis often requires larger sample sizes than simple regression analysis, especially when testing for higher-order interactions or when there are multiple moderators.

5. Multicollinearity: If the moderator variable is highly correlated with the predictor variable, this can lead to multicollinearity, which can make it difficult to estimate the moderation effect accurately.
6. Model complexity: Adding interaction terms to a statistical model can increase its complexity, which can make it more difficult to interpret and can lead to overfitting.

Despite these limitations, moderation analysis can provide valuable insights into the conditions under which a relationship between two variables is stronger or weaker, and can help researchers develop more effective interventions or treatments.

3 Mediation

Mediation refers to a causal process in which the relationship between two variables is explained by a **third** variable, called a mediator. Mediation analysis is used to test for the presence and strength of this causal process.

To test for mediation in a statistical model, researchers typically use a mediation model, which includes three variables: the independent variable (X), the mediator variable (M), and the dependent variable (Y). The mediator variable is assumed to be affected by the independent variable, and in turn, affects the dependent variable. The strength of the mediation effect is typically estimated using regression analysis.

Mediation analysis is important in statistics because it allows researchers to identify the mechanisms through which an independent variable affects a dependent variable. This can provide insights into the underlying processes that drive a particular phenomenon and can inform the development of interventions or treatments that target specific mediators.

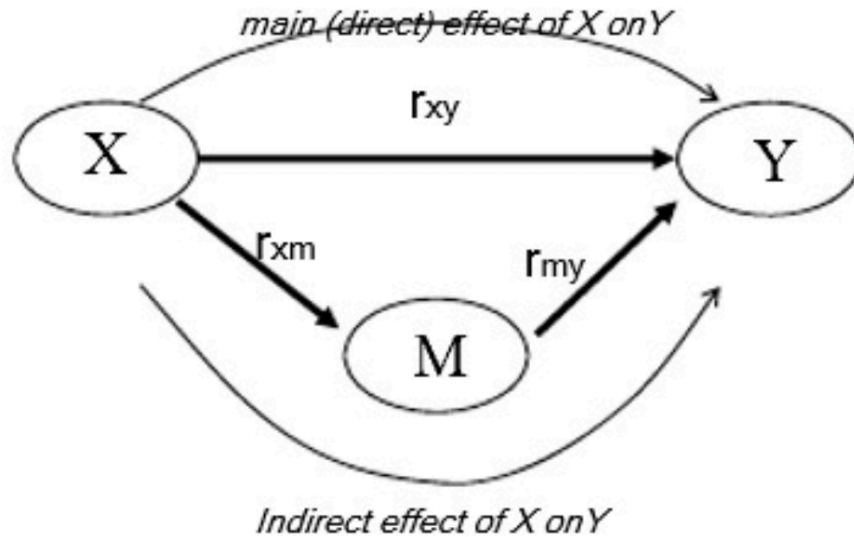


图 2: Moderation

3.1 Fit Mediation Model

The data in “mediation_data.sav” has information about customers’ satisfaction, relationship with customers and discounts offered. We are interested in whether discount can mediate the influence of customer relationship on their satisfaction about some consumer goods.

```
currentData <- bruceR::import("mediation_data.sav")
attach(currentData)
```

Use Baron and Kenny method to conduct mediation analysis. Three steps are:

1. IV could predict DV.
2. IV could predict Mediator.
3. Mediator could predict DV on the presence of IV.

Besides, the predictive power of IV is then reduced,

- If IV is no longer a significant predictor → full mediation.

- If IV is still significant but reduced in magnitude → partial mediation.

Note that the effect of mediator on DV should primarily be based on IV!

3.1.1 Direct Effect

Check without any consideration of mediation as a baseline model, which is the relationship between Y and X.

```
model_0 <- lm(Satisfaction ~ Relationship)
bruceR::GLM_summary(model_0)
```

3.1.2 Internal Correlation

```
# M as a function of X
model_M <- lm(Discount ~ Relationship)
bruceR::GLM_summary(model_M)
```

Significant result here is the necessity for further mediation analysis. If not significant, no need to move forward.

3.1.3 Full Model

```
# Y as a function of both X and M
model_Y <- lm(Satisfaction ~ Relationship + Discount)
bruceR::GLM_summary(model_Y)
```

Here, it's a case of partial mediation. And the adjusted R^2 is greatly improved, indicating that the mediation model is far more effective than the baseline model.

3.2 Mediation Effect

Mediation analysis is based on both the full model and the internal-correlation model. Here are more important indicators measuring the mediation effect:

- ACME, average causal mediation effect, the indirect effect of X on Y.
- ADE: average direct effect.
- Total Effect: the total effect of X on Y, the sum of direct and indirect effect.
- Proportion Mediated: how much percentage of indirect effect within the total effect.

```
library(mediation)
mediation_results <- mediate(model_M, model_Y, treat = "Relationship", mediator = "Disc
  boot = TRUE, sims = 500)
summary(mediation_results)
```

3.3 Appendix for Mediation

3.3.1 Mediation & Causality

Mediation analysis can reveal the presence of a causal relationship between two variables. A causal relationship exists when a change in one variable (the independent variable) leads to a change in another variable (the dependent variable). Mediation analysis examines the causal mechanism that links the independent variable to the dependent variable through one or more intervening variables (the mediators).

To establish a causal relationship, mediation analysis must satisfy three conditions:

1. The independent variable (X) must be significantly related to the mediator variable (M).

2. The mediator variable (M) must be significantly related to the dependent variable (Y), after controlling for the effect of the independent variable (X).
3. The relationship between the independent variable (X) and the dependent variable (Y) must be significantly reduced or eliminated when the mediator variable (M) is included in the model.

If these conditions are met, it provides evidence that the independent variable causes changes in the mediator variable, which in turn causes changes in the dependent variable. This suggests the presence of a causal relationship between the independent and dependent variables, mediated by the mediator variable.

It is important to note, however, that mediation analysis alone cannot establish causality with certainty. Other factors, such as confounding variables or reverse causality, may also influence the relationship between variables. Therefore, researchers should use caution when interpreting the results of mediation analysis and consider other forms of evidence, such as experimental designs. There are several limitations to using mediation analysis to establish causality:

1. *Directionality*: Mediation analysis assumes that the causal relationship flows from the independent variable to the mediator variable to the dependent variable. However, in some cases, the causal relationship may flow in the opposite direction or may be bidirectional, making it difficult to determine the direction of causality.
2. *Third variables*: Mediation analysis assumes that there are no third variables (confounding variables) that are responsible for the observed relationship between the independent variable and the dependent variable. If there are confounding variables that have not been accounted for, the observed mediation effect may be due to these variables rather than the mediator variable.
3. *Measurement error*: Mediation analysis assumes that all variables are

measured without error. However, in practice, measurement error is often present, which can lead to biased estimates of the mediation effect.

4. Sample size requirements: Mediation analysis often requires larger sample sizes than simple regression analysis, especially when testing for higher-order mediation or when there are multiple mediators.
5. Model complexity: Adding mediator variables to a statistical model can increase its complexity, which can make it more difficult to interpret and can lead to overfitting.
6. Nonlinearity: Mediation analysis assumes that the relationships between variables are linear. However, in some cases, the relationships may be nonlinear, which can lead to biased estimates of the mediation effect.

Despite these limitations, mediation analysis can provide valuable insights into the causal mechanisms that underlie relationships between variables. However, researchers should use caution when interpreting the results and consider other forms of evidence, such as experimental designs, to establish causality more convincingly.

3.3.2 Sobel Test

It's outdated admittedly. In early periods, the statistic constructed in Sobel Test was considered to follow a normal distribution, but that was proved wrong later. Hence, in Sobel Test, the probability of committing type I error increases before you even know it.

Sobel Test uses z -statistic:

$$z = \frac{\hat{a}\hat{b}}{\sqrt{\hat{a}^2 s_b^2 + \hat{b}^2 s_a^2}}$$

The code is exhibited below but I don't wanna run it.

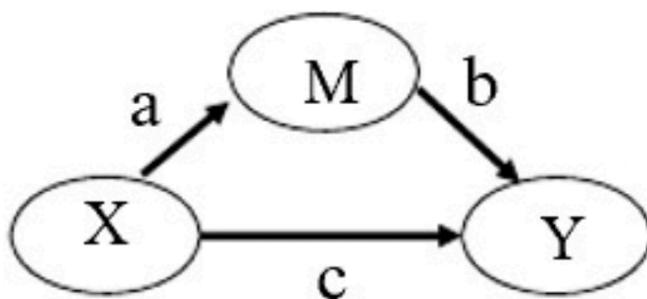


图 3: Sobel Test

```

library(bda)
mediation.test(Discount, Relationship, Satisfaction)
  
```

Lastly, don't forget to

```
detach(currentData)
```

4 容易踩的坑

调节与中介都容易栽的几件事

1. 把交互显著当作“调节存在”，忘了画 **simple slope**: b_3 显著只说明斜率随 Z 变化，但 ** 怎么变 ** 要看 simple slope。可能 Z 高时正向显著、 Z 低时不显著（增强效应），也可能两端方向相反（buffering / crossover）。结论完全不一样。
2. 中介不做 **bootstrap CI**: 很多老教材停留在 Baron & Kenny 的「四步法 + Sobel test」，但近 15 年的方法论文献一致推荐 bootstrap CI 检验间接效应。10000 次重抽 + 95% percentile CI 是标配。
3. 把中介和调节简单堆叠成“调节中介”模型不查可识别性: moderated mediation ($X \rightarrow M \rightarrow Y$ 但 a 或 b 路径被 Z 调节) 需要明确指定哪一条路径被调节、 Z 在哪个方程里出现。盲目把所有变量塞 PROCESS 模板里跑不一定能识别。建议先画路径图，确

认每条箭头都有理论支持再上模型。

4. **中介的因果方向只来自理论**： $X \rightarrow M \rightarrow Y$ 跟 $X \rightarrow Y$ (M 是结果的一部分) 在横截面数据下完全等价——区分只能靠时间顺序 (纵向设计) 或随机干预 (实验)。横截面中介系数显著 **不能** 直接得出因果结论。