

2023-1 StatisticFinalTerm Note

Version 1

Jm Su

<https://hello.jmsu.me>

June 2023

Contents

1	t-检验 (t-검증)	3
1.1	t-检验方法	3
1.2	t 与 z 的差异	3
1.3	关于单双侧检验	4
2	卡方检验 (x^2-검증)	4
2.1	卡方检验方法	4
3	ANOVA 분석 (F-statistic)	5
3.1	假设陈述	5
3.1.1	什么是处理效果 (Treatment effect) ?	6
3.1.2	组内, 组间均方	6
3.2	F-检验 (F-검증)	7
3.2.1	通过组内, 组间均方度量处理效果和随机误差	7
3.2.2	F 值的计算	7
3.3	ANOVA 的不足	8
3.3.1	评估处理效应的大小: η^2	8
3.3.2	事后比较 (사후비교)(post hoc comparisons)	8
4	Linear Regression	8
4.1	线性回归模型	9
4.2	计算线性回归常数	9
4.3	验证总体	10
4.3.1	线性回归和 ANOVA	10
4.3.2	R^2 决定系数	11
4.3.3	t-检验 (检验回归系数 b_{yx} 的显著程度)	11

5	Logistic Regression	11
5.1	逻辑回归模型	11
5.2	Odd ratio	12
5.3	对回归系数 b 的解释	12
6	附表	13

1 t-检验 (t-검증)

t 检验的对应变量类型应该是X: 定类变量 对Y: 连续型变量, 因为其验证的是两个集团间的平均的差异

1.1 t-检验方法

t-수치

$$t = \frac{\bar{Y} - \mu_y}{\sigma_{\bar{y}}} = \frac{\bar{Y} - \mu_y}{s_y / \sqrt{N}}$$

\bar{Y} : 표본평균 (样本平均)

μ_y : 모집단평균 (总体平均) (期望)

$\sigma_{\bar{y}}$: 표준오차 (标准误差) 指每次抽样均值与总体均值 (期望) 间的差异 (的标准化)

s_y : 표본표준차 (样本标准差) 指每个抽样观测值与样本均值 间的差异 (的标准化)

N : 표본수 (样本数量)

t 检验用于检测单一平均或两集团平均差异

检验原理是, 通过 t 数值算样本与总体间的 (用标准误差衡量) 偏差是多少, 将 t 值带入 t 分布与 $t_{c.v.}$ (拒绝值) 进行比较, $t_{c.v.}$ 依据三个要素决定:

$t_{c.v.}$ 의 영향요소

1. 자유도 (df) (自由度)
2. 단측 or 양측검증 (单, 双侧检验)
3. α 수준

- 自由度 (df): 指得是样本的自由程度, 在 t 检验中 df 为 N-1, 因为最后一个决定的样本被之前的 N-1 个样本确定了他可能的取值范围, 所以自由度为 N-1.

- 单双侧检验: 根据假设的类型, $t_{c.v.}$ 在分布中的位置

- α 水准: 即人为设定的弃真概率 (Type 1 Error)

最后将 t-수치与 $t_{c.v.}$ 比较后, 得出是否拒绝 H_0 假设.

1.2 t 与 z 的差异

功能上, z 只可用于单一平均检验, t 可用于单一平均检验和两集团平均检验.

z-수치

$$z = \frac{\bar{Y} - \mu_y}{\sigma_{\bar{y}}} = \frac{\bar{Y} - \mu_y}{\sigma_y / \sqrt{N}}$$

计算方法上, z 和 t 唯一的区别是对于 $\sigma_{\bar{y}}$ 标准误差计算上的差异, z 利用总体 (모집단) 的 (σ_y) 标准差计算标准误差, 但是总体标准差一般情况下很难得到, 因此 t 采用了 (s_y) 样本标准差替代总体标准差.

与直觉相反的是, t 相比于 z 的优势在对于**小样本**的检验上依旧存在, 因为根据中心极限定理, 小样本的情况, 样本抽样很可能不符合正态分布, 因此 z 无法使用, 并且 t 因为使用了样本标准差作为计算值, 所以反而受抽样偏差影响更小.

在大样本的情况下, t 也具有相对 z 的优势, 即可以不考虑总体是否符合正态分布, 因为抽样自动实现了正态分布 (中心极限定理).

1.3 关于单双侧检验

单侧检验指拒绝值在概率分布的一侧, 即假设为 $H_0 < ?$ 或 $H_0 > ?$ 的情况

双侧检验指拒绝值在概率分布的两侧, 即假设为 $H_0 = ?$ 或 $H_0 \neq ?$

2 卡方检验 (x^2 -검증)

x^2 检验 (Chi-square) 的对应变量类型应该是X: 定类变量 对Y: 定类变量, 其通过判断集团观测频数与其期望频数间的差异而验证两个或多个集团间的关联性.

2.1 卡方检验方法

如上所述, 卡方检验通过判断频数 (빈도, Frequencies) 差异来确定集团间的关联性, 所以 x^2 -수치의 计算需要先输出变量的交叉表 (교차분석표, Cross-tabulation), 交叉表会显示出各变量的**观测频数** (f_{ij}), 并在表行 or 列的边界显示出行计和列计 (Marginals). 根据行计和列计, 我们可以计算出**期望频数** (\hat{f}_{ij}), 通过观测频数和期望频数完成对 x^2 的计算, 如下:

기대빈도

$$\hat{f}_{ij} = \frac{(f_{i\cdot})(f_{\cdot j})}{N}$$

f_{ij} : 관찰빈도 (Observed Cell Frequencies)

$f_{i\cdot}$: 观测频度的行计

$f_{\cdot j}$: 观测频度的列计

\hat{f}_{ij} : 기대빈도 (Expected Cell Frequencies)

x^2 -수치

$$x^2 = \sum_{i=1}^R \sum_{j=1}^C \frac{(\hat{f}_{ij} - f_{ij})^2}{\hat{f}_{ij}}$$

R: Row 指行

C: Col 指列

对于公式, 其加和了各个格子 (房)(Cell) 的观测频数与期待频数间的差, 并通过自身平方并除以期待频数去掉了 \pm 号并标准化, 这意味着 x^2 统计量越大, 观测频数与期待频数间的差异也就越大.

换一种方式来说, 在自由度 (df), α 水准相同的情况下, x^2 越大, 越有可能拒绝 H_0 假设.

자유도 (df)

$$df = (R - 1)(C - 1)$$

注. 对于 x^2 分布, $N \geq 50$ 展现其分布的稳定性, 因此 x^2 检验最好适用于大样本

3 ANOVA 분석 (F-statistic)

ANOVA 区别与 t-检验和 z-检验, 它可以实现两个及以上集团间的平均差异比较, 变量类型为 X: 定类变量 Y: 连续型变量, 假设如下

3.1 假设陈述

对于 ANOVA 分析, 其包含了不同的解释方法和检验思路

ANOVA 가설진술방식 1

$$H_0 : \mu_j = \mu \text{ for all } j$$
$$H_1 : \mu_j \neq \mu \text{ for some } j$$

μ_j : 集团 j 的平均

以对比**平均差异**的方式建立假设，即零假设为：所有集团的平均相同

ANOVA 가설진술방식 2

$$H_0 : a_j = 0 \text{ for all } j$$
$$H_1 : a_j \neq 0 \text{ for some } j$$

a_j : 集团 j 的处理效果 (Population treatment effect)

以对比**处理效果**的方式建立假设

3.1.1 什么是处理效果 (Treatment effect) ?

(集团) 处理效果：集团平均与总平均的差异

$$a_j = \bar{Y}_j - \bar{Y}$$

一元变量的一般分析模型为

$$Y_{ij} = \mu + a_j + \varepsilon_{ij}$$

即，集团内的个体值由 总平均 + (集团) 处理效果 + 随机误差 决定

将以上公式合并，得出的一般模型为

$$Y_{ij} = \bar{Y} + (\bar{Y}_j - \bar{Y}) + \varepsilon_{ij}$$

3.1.2 组内, 组间均方

ANOVA 实际是通过 F-检验确定组内, 组间均方的比值确定处理效果的显著与否。三种假设虽然各有不同，但是只是解释视角的差异，本假设是最贴近 ANOVA 的实际计算方法的假设.

ANOVA 가설진술방식 3

$$H_0 : MS_B \leq MS_W$$
$$H_1 : MS_B > MS_W$$

MS_B : 组间均方
 MS_W : 组内均方

以对比组间, 组内均方的方式建立假设

3.2 F-检验 (F-검증)

3.2.1 通过组内, 组间均方度量处理效果和随机误差

所谓均方, 指平方和 (자승합)(Sum of squares)(SS) 的平均, 所以均方的计算基于平方和的计算, 如下

자승합 SS

$$SS_{Total} = \sum_{j=1}^J \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y})^2$$
$$SS_W = \sum_{j=1}^J \sum_{i=1}^{n_j} (Y_{ij} - \bar{Y}_j)^2$$
$$SS_B = \sum_{j=1}^J \sum_{i=1}^{n_j} (\bar{Y}_j - \bar{Y})^2$$
$$SS_{Total} = SS_B + SS_W$$

SS_T : 总平方和
 SS_W : 组内平方和 SS_B : 组间平方和

为何说平方和是处理效果和随机误差的度量呢, 公式中能看到组间平方和是通过计算各组平均和总体平均间的差异, 并将其平方加和. 同处理效果一样, 表示的是组间的差异性.

而如果对上节的一元一般变量模型公式进行左右互换, 消除可得到: 随机误差(ε_{ij})的计算方法为 $Y_{ij} - \bar{Y}_j$ 同组内平方和一样.

集团间的差异显著性即, 处理效果大于随机误差 (同假设.2), 而具体实现是通过比较组间平均和组内平均实现的, 如果组间平方和 (均方) 显著大于组内平方和 (均方), 那么可以判定: 集团间是存在显著变异性的, 即平均存在显著差异 (同假设.1).

3.2.2 F 值的计算

具体通过 F-分布检验显著性, 如下

F-검증

$$MS_B = SS_B/df = SS_B/J - 1$$

$$MS_W = SS_W/df = SS_W/N - J$$

$$F = MS_B/MS_W$$

MS 为均方, 在 3.12 节已经定义, F 值即组间均方和组内均方的比值, 也就是说组间均方越大, 组内均方越小的情况下, F 值越大, 集团间的显著性也越大.

3.3 ANOVA 的不足

如果仔细观察 3.1 节的三种假设, 可以发现 ANOVA 并不能测定

- 哪些集团间存在差异
- 处理效应大小的评估

因此需要其他统计值和事后比较来补足

3.3.1 评估处理效应的大小: η^2

$$\eta^2$$

$$\eta^2 = \frac{SS_B}{SS_T} = 1 - \frac{SS_W}{SS_T}$$

取值范围为 0~1, 指自变量对因变量处理效应的百分比, 0 即自变量对因变量没有影响 (完全独立).
一元 ANOVA 的情况下, 取值为 0.2(20%) 已经是不错的结果.

3.3.2 事后比较 (사후비교)(post hoc comparisons)

通过事后比较确定哪些集团间存在差异.
(사회통계학 -김상옥 p.142)

4 Linear Regression

线性回归假设自变量 X 到因变量 Y 的因果线性关系, 变量类型为 X: 连续型变量 Y: 连续型变量. 一元线性回归适用于单一 X, 多元线性回归适用于多 X (多自变量).

4.1 线性回归模型

관찰값 (Y_i) 모델

$$Y_i = a + b_{yx}X_i + e_i$$

Y_i : 관찰값 (观测值)

a : 절편 (截距)

b_{yx} : 회귀계수 (回归系数)

e_i : 무작위오차 (随机误差)(Random errors)(residual[残差])

X_i : 自变量取值

观测值指实际观测到的 Y 的取值; **截距**是线性回归方程中的常数之一, 指方程距离原点 0 的偏离;

回归系数是线性回归方程中的常数, 指自变量 X 对因变量 Y 的影响权重, 线性回归的核心就是估计出此权重, 完成线性方程; **随机误差**在线性回归中也叫**残差**, 指预测值与观测值间的距离 (差异).

예측값 (\hat{Y}_i) 모델

$$\hat{Y}_i = a + b_{yx}X_i$$

\hat{Y}_i : 线性模型预测值

与上面的观测模型比较, 可以看出预测模型去掉了随机误差, 因此也可得出随机误差等于

$$e_i = Y_i - \hat{Y}_i$$

这也很好解释, 因为真实观测值-模型预测值, 即误差, 也叫残差.

4.2 计算线性回归常数

회귀계수 b_{yx}

$$b_{yx} = \frac{s_{xy}}{s_x^2}$$

s_{xy} : xy 의 공변량 (协方差)

s_x^2 : x 의 변량 (方差)

절편 a

$$a = \bar{Y} - b_{yx} \cdot \bar{X}$$

通过以上公式计算出的回归直线叫做最小二乘回归线，其代表着这条直线拥有最小的随机误差 (e)(残差)，即公式 1 最小化.

$$\sum_{i=1}^n e_i^2 \quad (1)$$

4.3 验证总体

모집단모델
$Y_i = a + \beta_{yx}X_i + e_i$
b_{yx} : 회귀계수 (样本回归系数) β_{yx} : 모집단 회귀계수 (总体的回归系数)

上述计算的所有回归统计量都是基于样本变量 计算得出的，而其对于总体的显著性如何？这个问题需要通过检验过程得出.

4.3.1 线性回归和 ANOVA

线性回归和 ANOVA 有很大的相似，ANOVA 中的总平方和 (SS_T) 等于组内平方和 (SS_W) 加上组间平方和 (SS_B)，其表示了因变量 (Y_i)，因处理效果 (a) 加上随机误差 (ϵ) 对于其平均的总变异性.

对于线性回归模型，其因变量类型为连续型，因此其统计名称与 ANOVA 的分组平方和不同，主要使用如下统计名称

자승합 (SS) in Linear Regression
$SS_T = SS_{Error} + SS_{Regression}$ $SS_{Error} = \sum_{i=1}^n (Y_i - \hat{Y}_i)$ $SS_{Regression} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y}_i)$
SS_{Error} : 指预测值与观测值间的差异平方和，即残差平方和，简称 SS_E $SS_{Regression}$: 指预测值与观测平均值间的差异平方和，简称 SS_R

可以看出， SS_E 和 SS_R 等同于 ANOVA 的 SS_W 和 SS_B ，即随机误差和处理效果.

4.3.2 R^2 决定系数

在 ANOVA 中, η^2 用于测量自变量 X 对因变量 Y 的处理效应百分比, 其使用 SS_B/SS_T 计算得出, 即组间平方和占全体平方和的比例.

线性回归中, 计算处理效果的方法与 ANOVA 相同, 只要求出回归平方和占全体平方和的比例就可以了, 即:

R^2 결정계수 (Coefficient of Determination)

$$R^2_{yx} = \frac{SS_R}{SS_T}$$

R^2 的取值范围为 0~1, 表示因变量的总变异中可以被回归模型解释的比例.

附. 对 R^2 开根号, 可以得出 X 与 Y 的相关系数 (r_{xy}), 其取值为-1~1, 可表示负相关或正相关关系.

4.3.3 t-检验 (检验回归系数 b_{yx} 的显著程度)

b_{yx} 에 대한 t-검증

$$t = \frac{b_{yx} - \beta_{yx}}{SE(b)}$$

零假设为总体回归参数为 0, 即 XY 无线性关系, 因此 $\beta_{yx} = 0$.
 $SE(b)$: b 的标准误差 (计算公式参照 p.158 공식 11-6)

5 Logistic Regression

逻辑回归作为线性回归的变种, 适用于 X: 定类或连续型变量 Y: 定类变量, 利用线性模型拟合因变量发生的概率.

5.1 逻辑回归模型

逻辑回归通过对线性模型进行 Sigmoid 变换得到可预测二分类概率的曲线模型, sigmoid 模型如下

Sigmoid model

$$S(x) = \frac{1}{1 + e^{-x}}$$

e : 自然常数

将线性模型传入 Sigmoid 函数进行转换, 即得到了 logistic 模型

Logistic model

$$\hat{Y} = \frac{1}{1 + e^{-(a+b_{yx}X_i)}}$$

\hat{Y} : 事件发生的概率, 取值范围 0~1

通过 logistic 模型, 可得出线性模型所代表的含义

logistic 模型的线性表示

$$\log\left(\frac{\hat{Y}}{1 - \hat{Y}}\right) = a + b_{yx}X_i$$

\hat{Y} 代表事件发生的概率, 那么与之相对 $1 - \hat{Y}$ 代表事件不发生的概率

5.2 Odd ratio

Odd ratio 指赔率比, 即事件发生概率和事件不发生概率间的比值, 其取值范围为 $0 \sim \infty$.

如上 logistics 模型线性表示公式能看出, 线性回归模型对应的结果正是对数后的赔率比, 即 $\log(\text{Odd ratio})$, 因此在解释自变量 X_i 对因变量 Y_i 的影响时, 其实是赔率比的变化.

5.3 对回归系数 b 的解释

在解释时, 一般对 $\exp(b)$ 进行解释, 因为其消除了 \log 函数, 因此 $\exp(b)$ 的值即赔率比 (Odd ratio). 可以将其解释为: 当自变量增大 1 个单位, 事件发生的概率是事件不发生概率的? 倍.

即当 $\exp(b) < 1$ 时, 自变量的增加与因变量发生概率间成负关系; $\exp(b) > 1$, 自变量的增加与因变量发生概率间成正关系; 当 $\exp(b) = 1$, 即自变量和因变量之间无关联.

6 附表

各统计方法适用类型表

X \ Y	定类型	连续型
	定类型	连续型
定类型	X^2 -검증 로직회귀	z,t-검증 ANOVA(F-검증)
连续型	로직회귀	선형회귀