

重点：线性回归

2022年3月13日 星期日 21:49

① $r = \frac{Cov(x,y)}{\sqrt{S_x^2 S_y^2}}$ (相关系数)

② $b_1 = r \frac{S_y}{S_x}$

③ R^2 : 决定系数, 越高越好
(回归关系可以解释因变量多少变异)

① 代码过程 (一元)

$\star \rightarrow cor(x,y)$ 看相关系数
 $fit \leftarrow lm(y \sim x, data =) \rightarrow$ 看 Adjusted R-squared 大小 (higher is better)

residuals (fit) 看残差

predict (fit, dataframe (x=c(1,2,3...)))

\star Summary (fit) 看结果

② 代码 (二元)

$fit2 \leftarrow lm(data, formula = y \sim x + I(x^2))$

③ 代码 (非线性)

① $fit3 \leftarrow nls(data, y \sim exp(r * x), start = list(r =), alg = "plinear")$

④ 转化为 log10

回归诊断

① residuals are normally distributed

plot (fit, 2) Q-Q图

② errors are independent

plot (fit, 1) \rightarrow 残差与预测值毫不相干

③ 同方差 homoscedasticity (同②)

残差不随预测值 \uparrow 而 \uparrow

④ hist (resid(model))

看残差正态性

correlation analysis

连续变量

两组均为正态

等方差

符合

不符 (non-parametric)

Pearson

Spearman

$cor(x,y, method = "pearson")$

$cor(x,y, method = "spearman")$

默认 method = 'pearson'

画图

① geom_smooth (method = 'lm',

formula = y ~ x,

se = F/T)

geom_smooth (method = 'gam',

formula = y ~ x + I(x^2),

se = F/T)

\star ggpmisc 中 stat_poly_eq 在图上标表达式

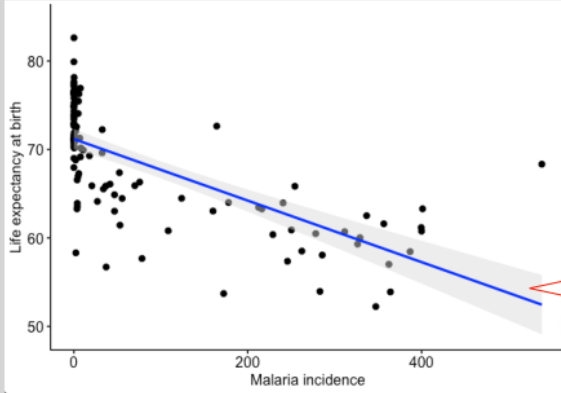
```
stat_poly_eq(aes(label = paste(..eq.label.., ..adj.rr.label.., sep = '-----'), color = "red", label.x = 0.25, geom = "label", formula = y ~ x + I(x^2), parse = T))
```

② abline (fit) 一元

\star lines (x, fitted (fit2)) 二元
(x轴) (y轴)

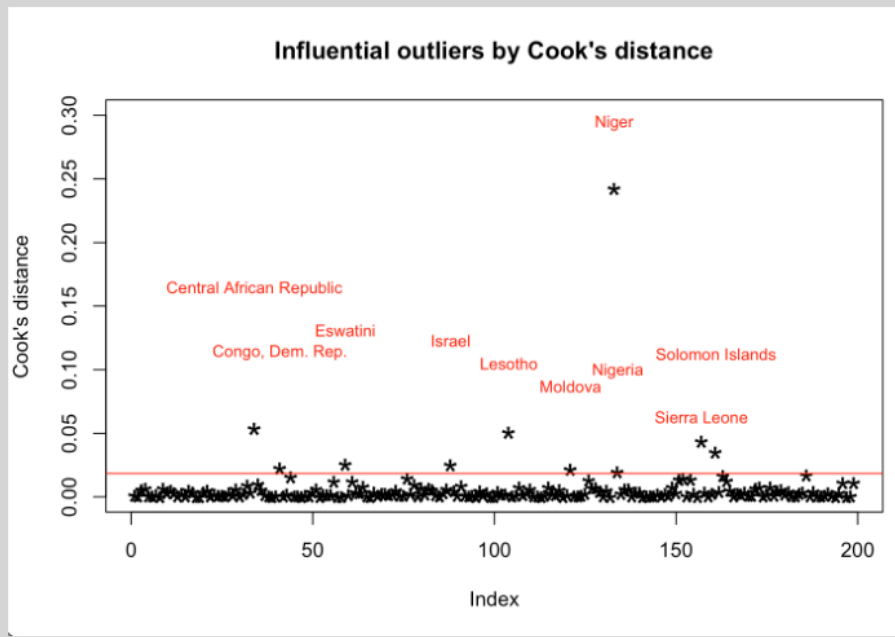
③ ggpubr 中 ggscatter

```
>ggpubr::ggscatter(WBd, x = "malaria_incidence", y = "lifeexpectancy", add = "reg.line", add.params = list(color = "blue", fill = "lightgray", conf.int = TRUE) + labs(x = "Malaria incidence", y = "Life expectancy at birth", colour = ""))
```



① 画 outliers 图

```
>cooks_d <- cooks.distance(lm(lifeexpectancy ~ fertility, WBd))
>sample_size <- nrow(WBd)
# Not NAs
>not_nas <- which(!is.na(WBd$lifeexpectancy) & !is.na(WBd$fertility))
# Plot Cook's distance
>plot(cooks_d, pch = "x", cex = 2, main = "Influential outliers by Cook's distance", ylim=c(0,1), ylab = "Cook's distance")
>abline(h = 4/sample_size, col = "red") # add cutoff line
>text(x = 1:length(cooks_d), y = cooks_d + 0.05, col = "red", cex = 0.8, labels = ifelse(cooks_d > 4 / sample_size, names(cooks_d), "")) # add labels
```



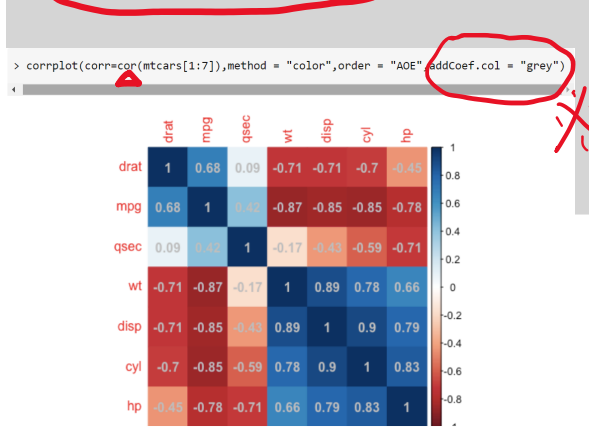
the graph is added to an existing plot

② corplot 包中 corplot 函数

```
>corplot::corplot(cor(WBd[withNAs == 0, 3:8]), order = "AOE", method = "ellipse", type = "upper", tl.pos = "d", add = TRUE, method = "number", type = "lower", diag = FALSE, tl.pos = "n", cl.pos = "n")
```

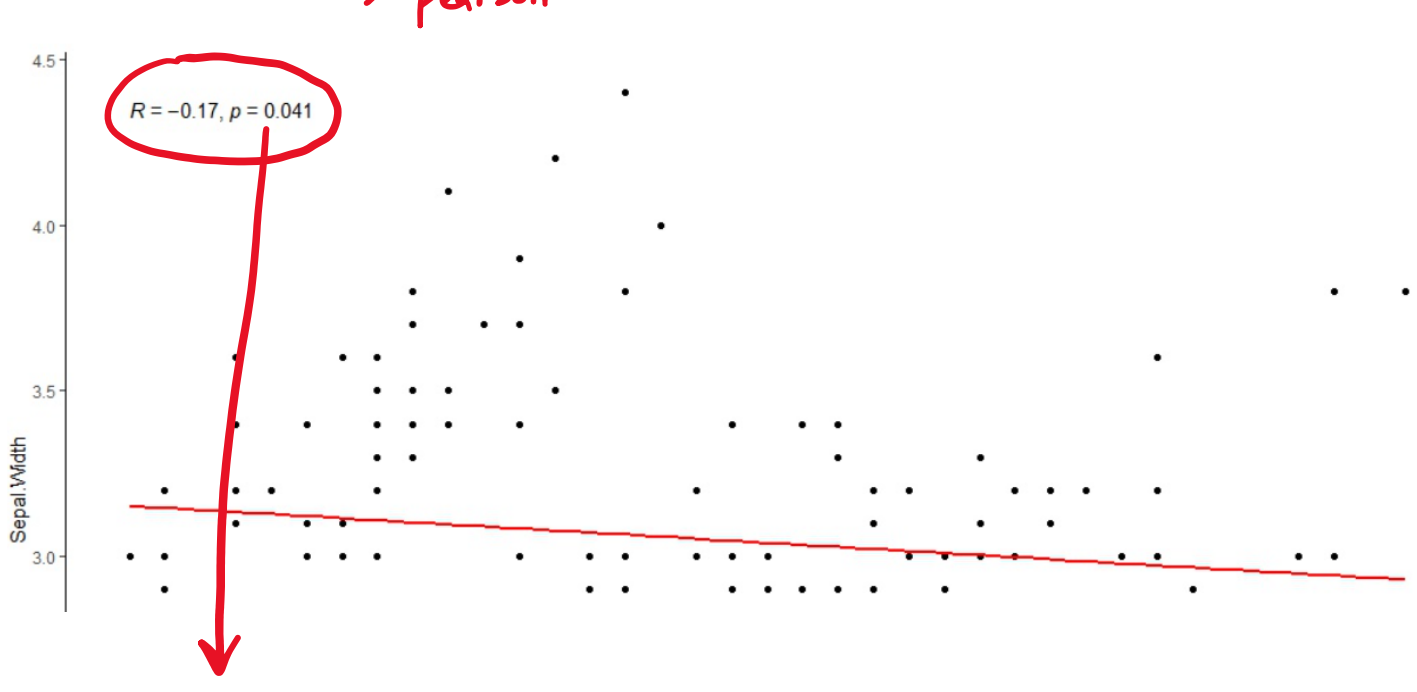
图例位置

method = "color"



③ 标相关系数 + p值

```
stat_cor(data = iris, method = "spearman")
```



讨论两变量是否相关必须讨论显著性水平, 不谈p值只谈相关系数大小是无意义的, 两者之间的相关关系可能只是偶然因素引起的, 所以我们要对两个变量之间的相关关系的显著性水平进行判断;

采用假设检验的方法:

原假设H0: R=0 两变量之间不存在线性关联

备择假设H1: R不等于0, 两变量之间存在线性关联