

دیت یلک؟

توضیح دهید

ا) غلط. به طریقی نمی توان در این باره نظر داد چون اگر dataset جدید outlier داشته باشد در این صورت خطای train می تواند بیشتر شود؛ پس ما این خطای train به مقادیر آن train می شود وابسته است.

b) دیت. طبق تعریف داریم رابطه ای داریم $e_i = y_i - x_i^T \hat{\beta}$

یا از مقادیر regression خطی از رابطه (*) تخمین زده می شود:

$$(*) \begin{cases} \hat{\beta} = (X^T X)^{-1} X^T Y \\ X = [x_i^T]_{N \times n} \rightarrow \text{(تعداد predictor ها)} \\ Y = [y_i]_{N \times 1} \rightarrow \text{تعداد داده ها} \end{cases}$$

(از راسته $X(X^T X)^{-1} X^T$)

$$\Rightarrow \sum_{i=1}^N y_i - x_i^T \hat{\beta} = \sum_{i=1}^N e_i = 0 \quad \checkmark$$

پس در این روش صیه می آید residual ها صفر است.

c) غلط. استاندارد سازی (features) فقط به مقدار regression در بین متغیرها که یک می باشد و اغلب آن بین از train می آید logistic regression

الزامی نیست.

d) دیت. $\rho = \text{correlation coefficient}$ نشان دهنده رابطه ای خطی میان دو متغیر است که در صورت همبستگی بین ارتباط خطی بین دو متغیر وجود ندارد اما ممکن است که بین دو متغیر (رابطه ای غیر خطی) وجود داشته باشد.

$$J = e^T e + \lambda \beta^T \beta \quad \text{ridge regression}$$

تابع هزینه در ridge regression

c) غلط. داریم:

اگر predictor ما نتواند خطای ss را کاهش دهد با افزایش λ regularization، مقادیر predictor کاهش یافته و تاثیرش در تخمین کم می شود. اگر λ را تا 0 هم درست است یعنی از ضرایب موجود در رابطه - به صفر نزدیک می شوند (صفری $\leftarrow 0$) اما - بطور کلی به صفر نمی روند.

f) $R^2 = 0 \leftarrow$ رابطه ای خطی بین دو متغیر وجود ندارد. \nleftrightarrow هیچ رابطه ای (و انش رابطه ای غیر خطی) بین دو متغیر وجود ندارد.

پس محتمل است رابطه ای دیگر (غیر خطی) بین دو متغیر وجود داشته باشد.

g) غلط. $R^2 = 0.711 \leftarrow |R| = 0.89$ مقدار رابطه ای خطی دارد که ما با R^2 مقایسه می توانیم بگویم رابطه ای خطی داریم یا نه. دیت. 0.88 دیت. 0.88

ادامه سوال ①

و بنابراین فقط می‌توانیم بدون توضیح به جهت بلوغ قدرت رابطی طی قدرت است بین وجود ترسب (راهله باطلی قسب)

با توجه به این موضوع دقت است.

مقدار مشاهده (observed)

$$e_i = y_i - \hat{y}_i$$

مقدار تخمین (regression value)

طبق تعریف برای اسدالها داریم:

به طبق تعریف اسدالها در صورت نامک مقادیر واقعی و تخمین از

درست ⑦

d) $y = A_0 + A_1 x$ (طبی گران داده)

evaluation score $\hat{y} = a_0 + a_1 x$

$$a_1 = \frac{S_y R}{S_x}$$

$$a_0 = \bar{y} - a_1 \bar{x}$$

$$\bar{x} = -0.0833$$

$$\bar{y} = 3.9983$$

$$a_1 = 0.1325$$

$$T = \frac{a_1 - 0}{\frac{S_y}{S_{x a_1}}} = \frac{0.1325}{0.032} = 4.141$$

	estimate	std. error	T-value	$P_r \{ > H1 \}$
intercept	4.010 ^a	0.025	157.21	0
adornment	0.1325	0.032	4.141	0

②

$$\begin{cases} H_0: \alpha_1 = 0 \\ H_1: \alpha_1 > 0 \end{cases}$$

$$T = 4.141$$

$$df = n - 2 = 461$$

ب) مطابق درین دایره:

test + یک طرفه

$$(\alpha = 0.05)$$

$$P_{value} = Pr\{t > 4.141\} < 0.05 \Rightarrow$$

یعنی فرض صفر رد می شود. $proof$ یعنی این فرض دلیل وجود دارد (رابطه بین $adornment$ و $evaluation$ مثبت است).

hypothesis test

$$95\% \rightarrow \alpha = 0.05 \quad t_{461}^* = 1.96$$

confidence interval

$$ME = \alpha_1 \pm 1.96 s_{\alpha_1} \rightarrow (0.06978, 0.1952)$$

کنش است.

a)

$$\hat{Y} = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \alpha_4 x_4 + \alpha_5 x_5 + \alpha_6 x_6$$

x_6 blood pressure (b.p)

	estimate	st error	t value	$Pr\{t > t \}$
(intercept)	-80.41	14.53	-5.6	0.0
x_1 ← Smoke	0.44	0.03	15.26	0.0
x_2 ← Exercise hours	-3.33	1.13	-2.95	0.0033
x_3 ← age	-0.01	0.09	-0.1	0.9170
x_4 ← height	1.15	0.21	5.43	0.0
x_5 ← weight	0.05	0.03	1.99	0.0471
x_6 ← water consumption	-8.40	0.95	-8.81	0.0

$$\hat{Y} = -80.41 + 0.44x_1 - 3.33x_2 - 0.01x_3 + 1.15x_4 + 0.05x_5 - 8.4x_6$$

b) P_{value} هر یک از متغیرهای (predictors) داریم:

for smoke $\rightarrow P_{value} \approx 0 \rightarrow$ یعنی این متغیر هم در معنی فارق است و اگرچه متغیرهای آب در نظر گرفته شود با افزایش هر یک از حجم smoke

$$t_{1229}^* = 1.962 \rightarrow$$

$$95\% \text{ میان } \alpha = 0.5$$

فارقون به طور متوسط 0.44 افزایش می یابد. (طبق محاسبات)

$$\alpha_1 \pm t_{1229}^* \cdot s_{\alpha_1} = (0.3811, 0.4989)$$

$$df = n - k - 1 = 1229$$

for water consumption $\rightarrow P_{value} \approx 0 \rightarrow$ این متغیر هم در معنی فارق است و اگرچه متغیرهای آب با افزایش این متغیرها

$$\alpha_6 \pm t_{1229}^* \cdot s_{\alpha_6} = (-10.264, -6.536)$$

$$\alpha = 0.5 \quad df = 1229$$

و حجم مصرف آب فارقون 8.40 کاهش می یابد. (مطابق محاسبات)

با افزایش باز می رود با رابطه

مثبت دارد.

$$R^2 = \frac{SSG}{SST} = 1 - \frac{SSR}{SST} = 1 - \frac{249.28}{332.57} = 1 - 0.75 = 0.25$$

$$R^2_{adj} = 1 - \left(\frac{SSR}{SST} \times \frac{n-1}{n-k-1} \right) = 1 - \left(0.75 \times \frac{1235}{1229} \right) \approx 0.247$$

③

اختلاف آلاینده در گروه ارقام مقایسه می‌گردد.
(goodr boadr medium)

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
class	2	1144	572		0.002326
Residual	425	39518	93		

b) $df_c = K - 1 = 3 - 1 = 2$

$df_f = n - 1 = 428 - 1 = 427 \rightarrow df_r = 427 - 2 = 425$

$MSG = \frac{SSG}{df_c} = 572 \rightarrow SSG = 1144$

$MSE = \frac{SSr}{df_r} = \frac{39518}{425} \approx 93$

بر ردی H_0 $P_r \{ F > 6.15 \} = 0.002326 < 0.01$
و H_1 در دست است یعنی دست‌کم بین دو گروه تفاوت دیده می‌شود.

c) $\alpha^* = \frac{\alpha}{K}$ (significance level)

دلیل کاهش خطای نوع دوم
 $\alpha^* = \frac{0.01}{3} \approx 0.003$

$df = 425 \rightarrow t_{df}^* = 2.95$

$(\bar{y}_1 - \bar{y}_2) \pm t_{df}^* SE = (4.926, 11.074)$

داریم: $SE = \sqrt{\frac{MSG}{n_1} + \frac{MSG}{n_2}}$

$H_0: \mu_1 = \mu_2 = 0$
 $H_1: \mu_1 - \mu_2 \neq 0$
 $SE \approx 1.042$

t-test بین دو گروه good و Medium

به جز با احتمال ۱٪ (۹۹٪) بازه مقدار برای مقایسه فرقی معنی‌دار
رد می‌شود بین دو گروه تفاوت داریم (۱۰۰٪)

اگر که مقادیر از این گروه معنی‌دار باشد، باید به خطای معمولی بینابین دقت کرد.

به مقادیر بدون این نیز باید دقت کرد. از آنجایی که داده‌ها به صورت random effect است، از random effect استفاده می‌گردد. از آنجایی که داده‌ها به صورت Multi-level Model است، از Multi-level Model استفاده می‌گردد.

random effect: فوایدی که predictor، توزیع نوال دارد.

(a) درستی نتایج به خصوصیات تست‌ها بستگی زیادی دارد. یکی از این خصوصیات که در تست‌ها به عنوان استقلال داده‌هاست، وجود

کوکیتهای وابستگی بین داده‌هاست. وجود این وابستگی تست‌ها را بی‌اعتبار می‌کند. داده‌های غیر مستقل، paired sample measures - repeated data

تست‌ها به طور دقیق‌تر به مقادیر وابسته می‌شوند. این وابستگی به صورت جفت‌شدگی، cross-classified و nested data از مقادیر اشتراکی پیچیده‌ای هستند که به استقلال وابسته

$$\tilde{\mu}_3 = \frac{\sum_{i=1}^N (x_i - \bar{x})^3}{N^3}$$

الف) معیار تحولی (skewness) ← معیار برای سنجش تقارن داده

معیار Kurtosis در واقع، heavy-tailed / light-tailed بودن داده سنجش
توزیع نرمال را نشان می‌دهد.

$$\tilde{\mu}_4 = \frac{\sum_{i=1}^N (x_i - \bar{x})^4}{N^4}$$

a) $\bar{x} = \frac{43 + 59 + \dots + 21}{8} = 34$, $s^2 = \frac{(43 - 34)^2 + \dots}{8}$ (7)

$$\frac{1498}{8} = 187.25 \rightarrow sd = s = \sqrt{187.25} = 13.6839$$

ب) نه. bootstrap، فقط موجود برای داده‌های کمی و کیفی است و این روش اطلاعات جدیدی ایجاد نمی‌کند، اگر در واقع داده‌ها
فرد (باینری) وجود داشته باشد (مانند داده‌های کیفی مانند bootstrap با کلاس نمی‌کند).

ج) با توجه به داده‌ها انتظار این است که به سمت راست متمایل باشد چون داده‌ها کوچک است و می‌تواند به طور تصادفی داده‌ها کوچکتر به نظر
رسند به علاوه این که bootstrap ویژگی‌های توزیع را حفظ نمی‌کند و توزیع داده اصلی را تقلیل به سمت (چوایی به راست) دارد پس
معیار bootstrap حفظ می‌شود.

د) زمان خرید نان (دیفیوژن و...) که از نظر توزیع نیز می‌توان آن تقریباً مانند توزیع uniform است.

$$CI : \bar{x} \pm 1.46 \times 4.85 = 34 \pm 8.05 = (25.95, 42.05)$$

(e)