

سوال ۱:

۱) در ابتدا با استفاده از دستور `table2array` داده‌های مربوط به متغیرهای `weeks` و `tgrams` را به ترتیب در آرایه‌های `weeks_array` و `tgrams_array` ذخیره می‌کنیم و سپس با استفاده از دستورات `mean`، `median`، `var`، `min` و `max` مقادیر میانگین، میانه، واریانس، مینیمم و ماکزیمم هر یک از دو متغیر را بدست می‌آوریم.

۲) یکی از روش‌های آزمون فرض آماری در مورد میانگین جامعه، آزمون تی (T-test) یا آزمون تی استیودنت (T-student) است. معمولاً در سه حالت از آزمون تی برای قضاوت در مورد میانگین استفاده می‌شود.

۱. آزمون تی تک نمونه‌ای (One sample t test) آزمون در مورد برابری میانگین جامعه یا مقدار ثابت و معلوم.

۲. آزمون تی دو جامعه مستقل (Two independent sample t test) آزمون در مورد برابری میانگین دو جامعه مستقل.

۳. آزمون تی برای زوج متغیرها (Paired sample t test) آزمون برابر میانگین دو متغیر از یک جامعه.

در همه این آزمون‌ها فرض نرمال بودن جامعه یا جامعه‌ها وجود دارد. همچنین ثابت بودن واریانس نیز برای هر سه حالت آزمون از فرض‌های اولیه است.

آزمون تی تک نمونه‌ای (One Sample T Test)

فرض کنید در یک جامعه نرمال می‌خواهیم میانگین جامعه را با یک مقدار مشخص (که ممکن است تحلیل گر حدس زده است) مقایسه کنیم. اگر میانگین جامعه را μ و مقدار حدس زده شده را با μ_0 نشان دهیم، فرضیات یعنی «فرض صفر» (Null Hypothesis) و «فرض مقابل» یا «فرض مخالف» (Alternative Hypothesis) مربوط به آزمون تک نمونه‌ای به صورت زیر نوشته می‌شوند.

$$H_0: \mu = \mu_0, H_1: \mu \neq \mu_0$$

بنابراین فرض صفر را به صورت «میانگین جامعه با مقدار μ_0 برابر است» و فرض مقابل را به صورت «میانگین جامعه با مقدار μ_0 برابر نیست» می‌خوانیم.

آماره آزمون برای حالت تک نمونه‌ای به صورت زیر نوشته می‌شود:

$$T = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$$

در این جا، منظور از \bar{x} میانگین نمونه، μ_0 مقدار حدسی برای میانگین و s انحراف استاندارد نمونه است. همچنین n نیز تعداد نمونه را نشان می‌دهد.

آماره آزمون در این حالت دارای توزیع t با $n-1$ درجه آزادی است. یعنی داریم:

$$T \sim t(n-1)$$

براساس مقدار احتمال (p-value) و احتمال خطای نوع اول (α) می‌توان در مورد رد فرض صفر تصمیم گرفت.

آزمون تی دو نمونه‌ای مستقل (Two Independent Sample T Test)

در این حالت با دو جامعه مستقل مواجه هستیم و می‌خواهیم میانگین آن دو را با یکدیگر مقایسه کنیم. آزمون تی با فرضیات زیر صورت بگیرد.

$$H_0: \mu_A = \mu_B, H_1: \mu_A \neq \mu_B$$

آماره آزمون در این حالت به صورت زیر نوشته خواهد شد.

$$T = \frac{(\bar{x}_A - \bar{x}_B) - (\mu_A - \mu_B)}{\sqrt{\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}}}$$

این آماره نیز دارای توزیع t با درجه آزادی زیر است:

$$\text{degrees of freedom} = \frac{\left(\frac{s_A^2}{n_A} + \frac{s_B^2}{n_B}\right)^2}{\frac{s_A^4}{n_A^2(n_A-1)} + \frac{s_B^4}{n_B^2(n_B-1)}}$$

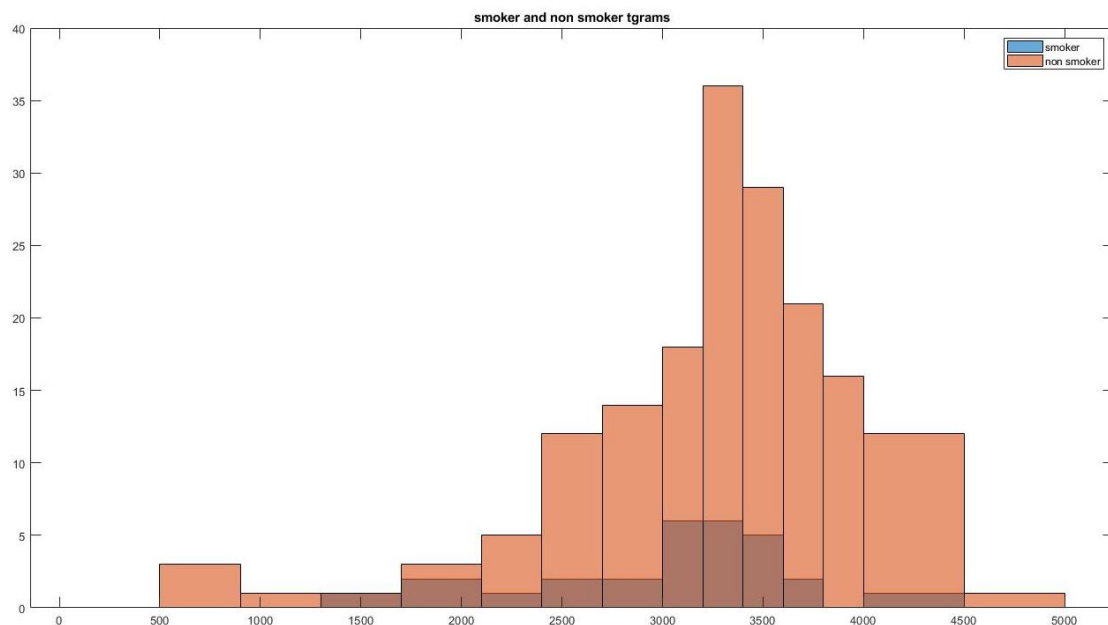
آزمون تی نمونه‌های وابسته (Paired Sample T Test)

آزمون تی دو نمونه‌ای یا آزمون تی زوجی از این جهت که برآورد واریانس در آماره آزمون آن متفاوت با حالت آزمون دو نمونه مستقل است مورد بحث قرار می‌گیرد. از این گونه آزمون بخصوص در زمانی که برای هر مشاهده دوبار اندازه‌گیری یک متغیر کمی صورت گرفته، استفاده می‌شود.

(۳ الف) در این قسمت فرض صفر را $\mu = 25$ و فرض مقابل را $\mu > 25$ در نظر می‌گیریم. با استفاده از مطالب گفته شده در قسمت ب، مقدار T و مرز ناحیه بحرانی را بدست می‌آوریم. از آن جایی که مقدار T درون ناحیه بحرانی است یا به طور معادل مقدار p_value از α کوچکتر است، فرض صفر رد می‌شود و در نتیجه گزاره صحت دارد.

(ب) در این قسمت فرض صفر را $\mu = 39$ و فرض مقابل را $\mu < 39$ در نظر می‌گیریم. با استفاده از مطالب گفته شده در قسمت ب، مقدار T و مرز ناحیه بحرانی را بدست می‌آوریم. از آن جایی که مقدار T درون ناحیه بحرانی است یا به طور معادل مقدار p_value از α کوچکتر است، فرض صفر رد می‌شود و در نتیجه گزاره صحت دارد.

(۴) در این قسمت فرض صفر را $\mu_0 - \mu_1 = 0$ و فرض مقابل را $\mu_0 - \mu_1 < 0$ در نظر می‌گیریم. با استفاده از مطالب گفته شده در قسمت ب، مقدار T و مرز ناحیه بحرانی را بدست می‌آوریم. از آن جایی که مقدار T خارج از ناحیه بحرانی است یا به طور معادل مقدار p_value از α بزرگتر است، فرض صفر تایید می‌شود و در نتیجه گزاره صحت ندارد.



سوال ۲:

(۱) ابتدا با استفاده از دستور `data, load` را در برنامه لود می‌کنیم. سپس با استفاده از دستور `table2array` داده‌های مربوط به متغیرهای مختلف را در آرایه‌هایی همنام با متغیرها ذخیره می‌کنیم. در متغیر `children` مقدار `NaN` را با مقدار ۴ عوض می‌کنیم. در ادامه با استفاده از دستور `unique`، تعداد مقدارهای مختلف هر متغیر را می‌یابیم. حال با استفاده از حلقه‌های `for`، به هر یک از مقادیر مختلف هر یک از متغیرها یک عدد طبیعی نسبت می‌دهیم. در انتها با کنار هم قرار دادن آرایه‌های مربوط به متغیرهای مختلف، ماتریس عددی `data` را می‌سازیم.

(۲) ابتدا با استفاده از دستور `randperm` برداری حاوی جایگاه هر نمونه می‌سازیم. سپس با استفاده از بردار بدست آمده، ماتریس `random_order_data` را از روی ماتریس `data_numeric_array` می‌سازیم. سپس داده اول ماتریس بدست آمده را به بخش آموزش و بقیه را به بخش تست اختصاص می‌دهیم.

(۳) ابتدا سطرهای مربوط به هر یک از مقادیر متغیر `class` و از روی آن احتمال هر یک از این مقادیر را بدست می‌آوریم. در ادامه برای هر یک از ویژگی‌ها یک ماتریس که تعداد سطر آن برابر تعداد مقادیر متغیر `class` و تعداد ستون آن برای تعداد مقادیر آن متغیر است، در نظر می‌گیریم و در خانه (i, j) آن، احتمال j امین مقدار آن متغیر به i امین مقدار متغیر `class` را قرار می‌دهیم.

(۴) ابتدا با استفاده از فرمول بدست آمده، احتمال هر یک از مقادیر متغیر `class` به شرط رخ دادن هر یک از نمونه‌های تست را بدست می‌آوریم. سپس ماکزیمم این مقادیر و شماره مربوط به مقداری از متغیر `class` که به ازای آن مقدار ماکزیمم بدست آمده است، برای هر یک از نمونه‌ها ذخیره می‌کنیم. این شماره‌ها همان مقادیر پیش‌بینی شده برای مقدار متغیر `class` هستند. در انتها با استفاده از دستور `plotconfusion`، `confusion matrix` مربوط به دسته‌بندی انجام شده را رسم می‌کنیم.

Confusion Matrix							
Output Class	1	992 33.5%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
	2	0 0.0%	874 29.5%	0 0.0%	103 3.5%	73 2.5%	83.2% 16.8%
	3	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	NaN% NaN%
	4	0 0.0%	101 3.4%	0 0.0%	813 27.5%	0 0.0%	88.9% 11.1%
	5	0 0.0%	0 0.0%	1 0.0%	0 0.0%	3 0.1%	75.0% 25.0%
Target Class							
1	2	3	4	5			

سوال ۳:

آ) ابتدا با استفاده از دستور `load, data` را در برنامه لود می‌کنیم. سپس با استفاده از دستور `table2array` داده‌های مربوط به متغیرهای مختلف را در آرایه‌هایی همنام با متغیرها ذخیره می‌کنیم. در ادامه با استفاده از دستور `unique`، تعداد مقادیر مختلف هر متغیر را می‌یابیم. حال با استفاده از حلقه‌های `for`، به هر یک از مقادیر مختلف هر یک از متغیرها غیر عددی یک عدد طبیعی نسبت می‌دهیم. در انتها با کنار هم قرار دادن آرایه‌های مربوط به متغیرهای مختلف، ماتریس عددی `data` را می‌سازیم.

ب) ابتدا با استفاده از دستور `randperm` برداری حاوی جایگاه هر نمونه می‌سازیم. سپس با استفاده از بردار بدست آمده، ماتریس `random_order_data` را از روی ماتریس `data_numeric_array` می‌سازیم. سپس `train` داده اول ماتریس بدست آمده را به بخش آموزش و بقیه را به بخش تست اختصاص می‌دهیم.

ج) با استفاده از دستور `fitlm` روش داده‌های آموزش (train) مدل خطی فیت می‌کنیم تا متغیر `area` را برحسب سایر متغیرها تخمین بزنیم. در اطلاعات خروجی این تابع در ستون `Estimate`، ضرایب فیت برای هر یک از ترم‌های مدل خطی، در ستون `SE` تخمین خطای آن ضریب، در ستون `t_stat` آماره `t` برای تست کردن این فرضیه که آن ضریب می‌توانست صفر باشد و در ستون `p_value` احتمال این که به طور تصادفی آن ضریب مقدار کنونی خود را داشته باشد، آمده است. بنابراین هر چه `p_value` بیشتر باشد، آن ویژگی کمتر در پیش‌بینی مساحت سوخته شده، تاثیر دارد. لذا در این دیتاست، به ترتیب داده‌های `DMC`، `X` و `month` بیشترین تاثیر را دارند.

د) ضریب تعیین (`R_squared`) مقدار متناسب تغییرات در متغیر جواب `y` را نشان می‌دهد که در مدل رگرسیون خطی به وسیله متغیر مستقل `X` بیان می‌شود. هر چه مقدار `R_squared` بیشتر باشد تنوع بیشتری با مدل رگرسیون خطی بیان می‌شود.

ه) با استفاده از `mdl.Fitted` و دیتا بخش آموزش دیتاست، با یک حلقه `for`، `MSE` مورد نظر را محاسبه می‌کنیم.

و) با استفاده از `mdl.predict` و دیتا بخش تست دیتاست، با یک حلقه `for`، `MSE` مورد نظر را محاسبه می‌کنیم.