

جلسه ۱ بیت و محاسبه آمار و احتمال صنفی

مهم ترین صنفی آمار و احتمال صنفی . مهم ترین تست های آماری روی بیان می کنیم .

χ^2 -Test

T-Test

تست کای ۲- (χ^2 Test)

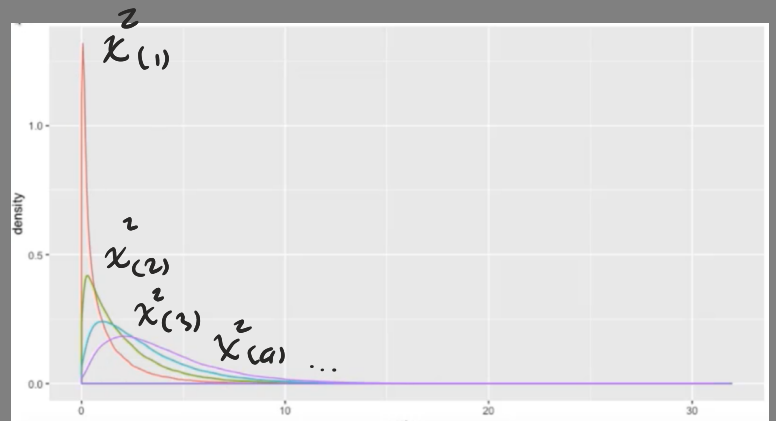
$X_1 \sim N(0,1) \Rightarrow X_1^2 \sim \chi^2(1)$ (degree of freedom) در بی آزادی

$X_2 \sim N(0,1) \Rightarrow X_1^2 + X_2^2 \sim \chi^2(2)$

\vdots
 $X_n \sim N(0,1) \Rightarrow X_1^2 + X_2^2 + \dots + X_n^2 \sim \chi^2(n)$

تعداد از یک به یک متصل باشند

مقادیر χ^2 مربعی است هر قدر صفت تغییر کند



اگر یک سری نمونه داشته باشیم که در چند دسته قرار بگیرند و هر نمونه در یک دسته قرار بگیرد.
 و کل نمونه ها m تا باشد.

$\underbrace{\quad}_{\text{بالا مثل } p_i \text{ این دسته}}$	$\underbrace{\quad}_{p_1}$	$\underbrace{\quad}_{p_2}$	\dots	$\underbrace{\quad}_{p_k}$
$m p_1$	$m p_2$	$m p_3$	\dots	$m p_k \sim \text{Expected}$
X_1	X_2	X_3	\dots	$X_k \sim \text{Observed}$

مساعد شده توسط نمونه گیری صدفی (واقعی)

آمار H_0 بجزر باشد. یعنی داده ها به همان احتمالات داده شده در هر دسته قرار گرفته باشند.

$$\sum_{i=1}^k \frac{(\text{observed}_i - \text{Expected}_i)^2}{\text{Expected}_i} \underset{m \rightarrow \infty}{\sim} \chi^2(k-1)$$

$$\sum_{i=1}^5 \frac{(X_i - m p_i)^2}{m p_i} \sim \chi^2(4)$$

بسی ۵: $k=5$

مثلاً اگر یک کالا توسط 3 کمپانی تولید شود و شخصی اظهار کند که α_1 درصد از برند 1، α_2 درصد از برند 2 و α_3 درصد از برند 3 خریدار آن کالاهای رویم و از جامعه نمونه گیری می‌کنیم. فرض را بر H_0 می‌گذاریم و به جای تست χ^2 می‌پردازیم. بعد از محاسبه χ^2 ، مقدار P-Value را می‌گیریم و P-Value است صحت صحت‌های شخصی را می‌سنجیم.

کاربرد اصلی χ^2

ما برای جداول وقوع 2×2 تست فیشر را داریم و برای جداول وقوع بزرگتر ما از تست کای-2 استفاده می‌کنیم. به نوعی کای-2 به‌عنوان تست فیشر است.

	α_1	α_2	α_3	
دارد χ دارد χ	0	0	0	برند 1
ندارد χ ندارد χ	0	0	0	برند 2
	0	0	0	برند 3

تست کای برای $m \rightarrow m$ مناسب است.
برای m های کوچک بهتر است که از تست فیشر استفاده کنیم.

تست فیشر
تست کای-2

حاشی درجه آزادی

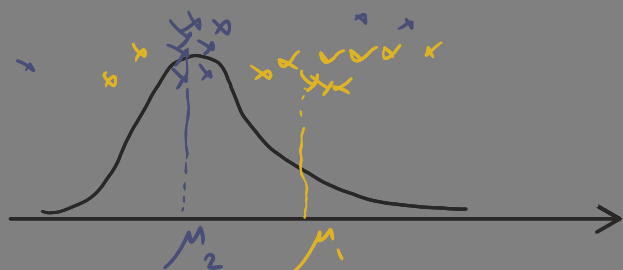
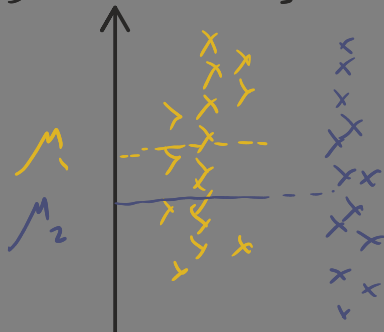
به جدول وقوع با اندازه‌های α و β درجه آزادی (K) از رابطه‌ی روبه‌رو بدست می‌آید.

$$K = (\alpha - 1)(\beta - 1)$$

خیلی وقت‌ها پرسش می‌آید که ما می‌خواهیم متوجه بشویم که آیا تفاوت‌ها معنادار هستند یا نه. برای مثال یا نایب عمرگاه‌ها می‌آید دو شرکت را مقایسه می‌کنیم.

تست -ی (t-test)

دو نمونه داریم (هر نمونه‌ای تولید می‌کند) (حاشیه با m). آیا این نمونه‌ها از توزیع‌های با میانگین یکسان آمده‌اند یا خیر؟



به سُرایی که بالا به مقادیر دیده می‌شود. ما روی یک توزیع یکسان در سیمبل جدا کرده ایم که میانگین‌های متفاوت دارند. دقیقاً همین جاست که اهمیت و کاربرد t -test برای ما مشخص می‌شود. با استفاده از تست t می‌توانیم در سیمبل رویه خوبی با هم مقایسه کنیم و به یکسان بودن توزیع همی ورود سیمبل پی ببریم.

همین اصلی طراحی یک آماره در این است که تحت برقراری H_0 ما توزیع sample باید ایم.

برای این که ما μ_1, μ_2 های نزدیک به هم داشته باشیم، خوشحالی کنیم و آن σ_1^2, σ_2^2 های کوچک داشته باشیم. در اینجا این sample نزدیک توزیع اومن چون پراکنده‌ی داده‌های او را هم کم هست.

$$\text{همی فرضی کلگه } H_0 \text{ برقراره} \Rightarrow \sigma_1^2 = \sigma_2^2 = \sigma^2$$

بریم سراغ توزیع

توزیع t

$$t = \frac{Z}{\sqrt{\frac{V}{n}}}, \quad \begin{matrix} Z \sim N(0,1) \\ V \sim \chi^2(n) \end{matrix}$$

n : اندازه‌ی نمونه

توزیع t

$$\Rightarrow t = \frac{Z}{S}$$

تنها پارامتره توزیع t به بی آزادی آن هست
 t (degree of freedom)

استدلالی واضح به عمل؟

آماره t

تک نمونه

یک سری داده داریم. فرض کنیم داده ما از توزیعی با میانگین μ_0 آمده اند.

$$t = \frac{\bar{X} - \mu_0}{\frac{S}{\sqrt{n}}}$$

تخمین ناآرایی از توزیع
جمعیت sample

$$\text{به بی آزادی} = n-1$$

توزیع t

معمولاً توزیع t خیلی شبیه به توزیع نرمال هست اما نه آنقدر شبیه به هم نیست که می‌تواند وقتی مرکزی که خیلی دور است.

اما به قدری که به بی آزادی نزدیک بشود و n بزرگ می‌شود در آن حالت توزیع t بسیار به توزیع نرمال نزدیک می‌شود.

دو نمونه

زمانی که می‌خواهیم در مورد داده‌ها مقایسه کنیم. H_0 : این دو نمونه از توزیع‌های با میانگین یکسان آمده‌اند.

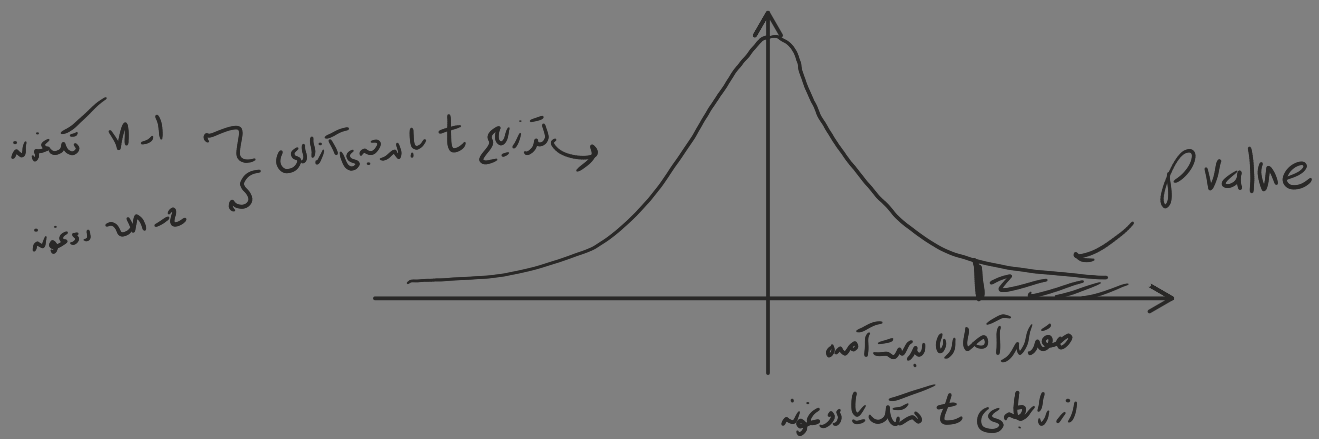
$$t = \frac{\bar{X}_1 - \bar{X}_2}{S_p \sqrt{\frac{2}{n}}}, \quad S_p = \sqrt{\frac{S_{X_1}^2 + S_{X_2}^2}{2}} \Rightarrow t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{S_{X_1}^2 + S_{X_2}^2}{n}}}$$

S^2 : واریانس sample های ما

درجه‌ی آزادی: $2n - 2$
توزیع t

حالا فرض کنیم با توجه به آماره داده‌ای که داریم رد نیست یا بپذیریم و بینیم که $P\text{-value}$ برای سنجش آماری این دو با توزیع t صفر است.

آنکه $P\text{-value}$ برای سنجش صفر آماره توزیع t به ازای‌های کافی که چید باشد H_0 واقعاً رد نیست.



اصلاً برای $P\text{-value}$ α برابر 0.05 است که به معنای 5٪ قابل قبولی است.

در زمان استاندارد لذت t نکات زیر باید رعایت شوند:
1. هرگز از بودن این نکات بی‌خبر نباشید.
2. استاندارد از دست t است.

1. میانگین‌های دو جامعه S باید توزیع نورمال داشته باشند.

2. وقتی به دست 2 نمونه استاندارد می‌کنیم فرضی برای این است که واریانس دو نمونه با هم برابرند. برای همین واریانس نام داده حل می‌دهد.

3. $\sqrt{2}$ و $\sqrt{2}$ در توزیع t باید لازم مستقل باشند.