

The background of the book cover consists of several overlapping, wavy, semi-transparent shapes in shades of light green, light blue, and light purple, creating a layered, abstract effect. The text is centered and rendered in a dark red, serif font.

Matemáticas e Estatística II

JOSÉ CARLOS DÍAZ RAMOS

Matemáticas e

Estatística II

JOSÉ CARLOS DÍAZ RAMOS

22 de febreiro do 2023

Este traballo ten licencia Creative Commons “Atribución-Non comercial-Compartir igual 4.0 Internacional (CC BY-NC-SA 4.0)”



Prefacio

Este traballo correspóndese cos apuntes da materia “Matemáticas e Estatística II” do Grao de Farmacia da Universidade de Santiago de Compostela. A versión máis actualizada pode atoparse en liña na dirección

<http://xtsunxet.usc.es/carlos/matematicas2/>

Estas notas están parcialmente baseadas nos apuntes da materia elaborados polos profesores Enrique Macías Virgós e Antonio Mariano Gómez Tato. O libro de referencia é [1]. Os exercicios están na súa meirande parte adaptados de [1] e de [2].

Índice xeral

1. Preliminares	1
1.1. Estatística descriptiva	1
1.1.1. Tipos de datos	1
1.1.2. Precisión	1
1.1.3. Medidas de tendencia central	2
1.1.4. Medidas de posición	3
1.1.5. Medidas de dispersión	3
1.1.6. Transformación de datos	4
1.2. Variables aleatorias	4
1.2.1. Distribución normal	5
2. Introducción á inferencia estadística. Estimación	7
2.1. Poboación e mostra	7
2.2. Estimación da media poboacional	9
2.2.1. Estimación puntual	9
2.2.2. Estimación por intervalos: coñecida a varianza poboacional	9
2.2.3. Estimación por intervalos: descoñecida a varianza poboacional	12
2.3. Estimación da varianza poboacional	14
2.3.1. Estimación puntual	14
2.3.2. Estimación por intervalos: coñecida a media poboacional	15
2.3.3. Estimación por intervalos: descoñecida a media poboacional	17
2.4. Estimación dunha proporción	18
2.4.1. Estimación puntual	18
2.4.2. Estimación por intervalos	19
2.5. Resumo de estimadores	21
3. Contraste de hipóteses	25
3.1. Contraste de hipóteses para a media da poboación	27
3.1.1. Contrastes bilaterais	27
3.1.2. Contrastes unilaterais	29
3.1.3. O valor P ou valor crítico	31
3.2. Contraste de hipóteses para a varianza	33
3.2.1. Contrastes bilaterais	33
3.2.2. Contrastes unilaterais	33
3.3. Contraste de hipóteses para unha proporción	34
3.3.1. Contrastes bilaterais	34

3.3.2. Contrastes unilaterais	34
3.4. Resumo de contrastes de hipóteses para unha poboación	36
4. Comparación de dúas poboacións	39
4.1. Comparación das medias de dúas poboacións con mostras independentes	39
4.1.1. Coñecidas as varianzas poboacionais	40
4.1.2. Descoñecidas as varianzas poboacionais pero supostas iguais	41
4.1.3. Descoñecidas as varianzas poboacionais	42
4.2. Comparación das varianzas de dúas poboacións con mostras independentes	45
4.3. Comparación de proporcións de dúas poboacións con mostras independentes	48
4.3.1. Intervalos de confianza	49
4.3.2. Contraste de hipóteses	49
4.4. Comparación da media con mostras emparelladas	51
4.5. Resumo de contrastes de hipóteses para dúas poboacións	53
5. A proba chi-cadrado	57
5.1. Contrastes de independencia para datos categóricos	57
5.2. Contrastes de homoxeneidade para datos categóricos	60
5.3. Bondade do axuste	63
6. Regresión e correlación	67
6.1. Regresión linear	68
6.1.1. Estimación dos valores	69
6.1.2. Covarianza e correlación	70
6.1.3. Regresión exponencial	73
6.1.4. Regresión potencial	75
6.2. Análise da varianza	75
6.2.1. ANOVA	76
6.2.2. Intervalos de estimación	80
7. Problemas para as clases interactivas	83
7.1. Intervalos de confianza	83
7.2. Contrastes de hipóteses	85
7.3. Contrastes de hipóteses para dúas poboacións	86
7.4. Problemas de repaso de estimación e contraste de hipóteses	88
7.5. Probas de homoxeneidade e independencia	89
7.6. Regresión linear e ANOVA	90
8. Exames resoltos	93
8.1. Exame 1	93
8.2. Exame de xuño de 2019	99
8.3. Exame de maio de 2021	105
9. Táboas estatísticas	113
Bibliografía	123

Capítulo 1

Preliminares

A Estatística é a rama das Matemáticas que estuda e interpreta os procesos aleatorios, para permitir deducir propiedades dunha poboación a partir dun subconxunto pequeno da mesma. Ademais de ter un corpo formal como parte das Matemáticas, a estatística é a miúdo empregada noutras ciencias co obxectivo de permitir establecer correlacións e dependencias entre diversos fenómenos físicos ou naturais.

1.1. Estatística descriptiva

A estatística descriptiva é a técnica matemática que organiza e describe un conxunto de datos co propósito de poder entendela con máis facilidade. A continuación presentamos algúns conceptos relevantes na estatística descriptiva.

1.1.1. Tipos de datos

Os datos poden ter diversa natureza:

- **Datos nominais**, que son etiquetas para distinguir a uns de outros, como as provincias de nacemento.
- **Escalas ordinais**, nas que se asigna unha orde, pero na que o número en si non ten relevancia, como a posición dun competidor nunha liga.
- **Escalas de intervalo**, que son medicións cuantitativas nas que se mide a diferenza entre dúas variables, como a temperatura en graos Celsius.
- **Escalas de razón**, que son escalas de intervalo cun cero absoluto, como a temperatura en graos Kelvin.

1.1.2. Precisión

A precisión entenderémola como o número de cifras representativas empregadas para expresar unha medida. Aínda que neste curso non faremos especial fincapé neste tema, convén ter en conta que os erros de precisión nos números se van propagando a medida que imos facendo operacións aritméticas. Para certos cálculos (p. ex. o coeficiente de

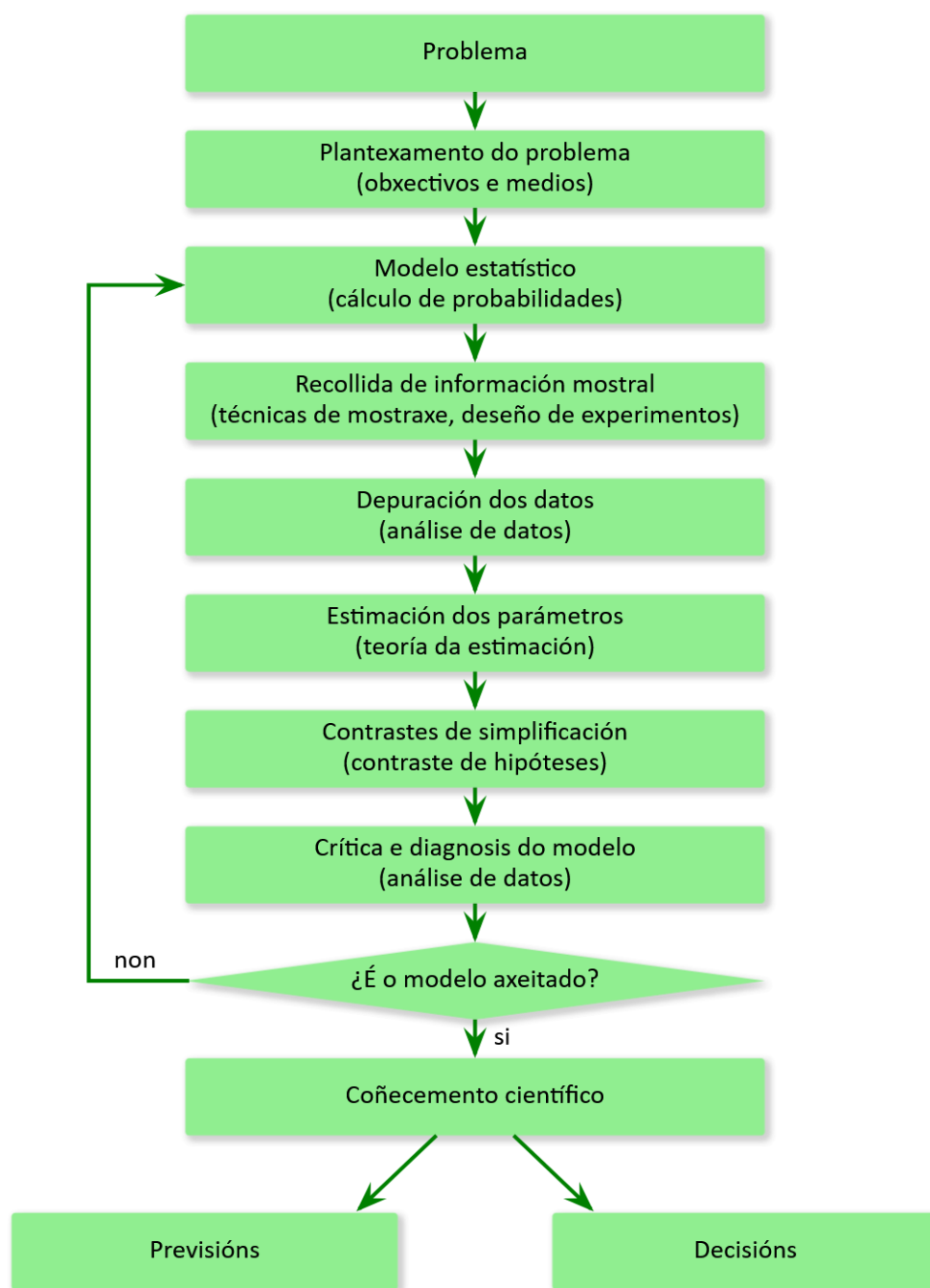


Figura 1.1: Diagrama de flujo do método estadístico

correlación (Subsección 6.1.2)) é necesario empregar unha cantidade suficiente de decimais para non chegar a resultados absurdos.

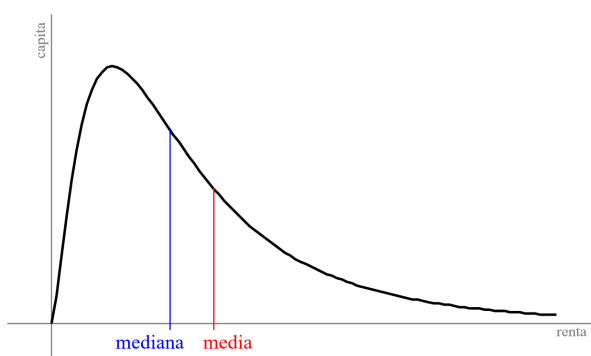
1.1.3. Medidas de tendencia central

- **Moda:** valor máis frecuente.
- **Mediana:** valor Md tal que, unha vez ordeados os datos, divide a estes pola metade.

- **Media:** é unha medida para datos obtidos como escalas de intervalo ou de razón, e que vén definida do seguinte xeito:

$$\bar{X} = \frac{X_1 + \dots + X_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Observación 1.1. A media e a mediana dan lugar a medidas similares en variables que se distribúen de xeito aproximadamente normal, como as alturas e os pesos dos seres vivos dunha determinada especie, ou os erros de medición. Para outro tipo de variables poden dar resultados moi distintos. Aínda que a media goza de máis popularidade e é sinxela de entender, hai ocasións en que a mediana resulta moito máis informativa e veraz.



Por exemplo, en termos económicos, a media soe dar información moi distinta á mediana. Unha medida habitual da economía é o produto interior bruto, ou a renda per capita. Esta última, que vén a ser unha media das rentas das persoas dun país, está sesgada cara ás élites dos ricos. Se por exemplo o 90% da xente perde poder adquisitivo, pero o 10% dos ricos se convirten en moito máis ricos, é perfectamente posible que a renda per capita aumente, dando impresión de que a economía mellora, a pesar de que o 90% da poboación lle vai peor. Non obstante, a mediana reflicte moito mellor a economía da maioría da xente, xa que nos dá a renda que divide á poboación en dúas metade do mesmo tamaño: a metade da poboación ten unha renda inferior a esa cifra, e a outra metade, superior. No caso anterior, a mediana da renda diminuíría, xa que a maior parte da xente perde poder adquisitivo. Con esta medida quedaría máis claro que é o que lle pasa á meirande parte da poboación.

1.1.4. Medidas de posición

- **Cuartís:** análogo á mediana, pero dividindo a distribución en cuartos Q_1 , Q_2 e Q_3 .
- **Percentís:** análogo á mediana e ós cuartís, pero dividindo a distribución en cen partes.

1.1.5. Medidas de dispersión

- **Rango:** diferenza entre o máximo e o mínimo.

- **Amplitude intercuartil:** diferencia entre Q_3 e Q_1 .
- **Desviación mediana:** mediana de $|X - Md|$.
- **Varianza:**

$$s_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}^2.$$

- **Desviación típica:** raíz cadrada da varianza.
- **Cuasi-varianza:**

$$s_{n-1}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{n}{n-1} s_n^2.$$

- **Cuasi-desviación típica:** raíz cadrada da cuasi-varianza.

Salvo que se especifique o contrario, neste curso asumiremos que s denota a cuasi-desviación típica, e s^2 a cuasi-varianza.

Proposición 1.2. *Séguese das fórmulas anteriores:*

- *A varianza, a cuasi-varianza, a desviación típica e a cuasi-desviación típica non poden ser negativas.*
- *Son cero se e só se tódolos datos son iguais á media.*

1.1.6. Transformación de datos

Se $Y = aX + b$ entón,

$$\begin{aligned}\bar{Y} &= a\bar{X} + b, \\ s_{n,Y}^2 &= a^2 s_{n,X}^2, \\ s_{n,Y} &= |a| s_{n,X}.\end{aligned}$$

A miúdo se empregan cambios para modifica-la media e a varianza:

- **Puntuacións desviadas:** $x = X - \bar{X}$. Así, $\bar{x} = 0$ e $s_{n,x} = s_{n,X}$.
- **Puntuacións tipificadas:** $z = \frac{1}{s_{n,X}}(X - \bar{X})$. Así $\bar{z} = 0$ e $s_{n,z} = 1$.

1.2. Variables aleatorias

Unha variable aleatoria pode describirse informalmente como unha variable que mide unha determinada característica numérica dunha poboación, de xeito que os seus valores dependen do resultado dun experimento aleatorio. Ó longo desta sección suporemos que X é unha *variable aleatoria absolutamente continua*, o que vén a querer dicir que dita variable toma os seus valores nun intervalo.

Toda variable aleatoria ten asociada unha **función de distribución** que vén dada por $F(x) = P(X \leq x)$, é dicir, $F(x)$ é a probabilidade de que a variable aleatoria X tome un valor menor ou igual ca x .

A **función de densidade** dunha variable aleatoria absolutamente continua é a derivada da súa función de distribución, $f(x) = F'(x)$.

A área baixo a gráfica da función da densidade nun determinado intervalo $[a, b]$ expresa a seguinte probabilidade:

$$P(a \leq X \leq b) = \int_a^b f(x)dx.$$

En particular, $P(X \leq b) = F(b) = \int_{-\infty}^b f(x)dx$.

Neste curso calcularanse moitas veces probabilidades do estilo

$$P(X \geq x) = \int_x^{\infty} f(x)dx.$$

Á función $P(X \geq x) = 1 - F(x)$ tamén se lle chama *función de supervivencia* da variable aleatoria X .

A media ou **esperanza** de X defínese como

$$E(X) = \int_{-\infty}^{+\infty} xf(x)dx.$$

A media ou esperanza da distribución denótase por μ .

A **varianza** de X defínese como

$$V(X) = \int_{-\infty}^{+\infty} (x - \mu)^2 f(x)dx = E(X^2) - E(X)^2.$$

A varianza dunha distribución denótase por σ^2 , e σ denotará a súa desviación típica.

O seguinte teorema dá unha idea de como se concentra a probabilidade dunha variable aleatoria arredor da media, sexa cal sexa a súa distribución.

Teorema 1.3. (Desigualdade de Chebyshev) Para unha variable aleatoria X satisfaise

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}.$$

Por exemplo, poñendo $k = 2$ na desigualdade de Chebyshev, obtemos $P(|X - \mu| \geq 2\sigma) \leq 1/4$, é dicir, que polo menos tres cuartas partes da probabilidade dunha variable aleatoria arredor da media están contidas entre $(\mu - 2\sigma, \mu + 2\sigma)$.

1.2.1. Distribución normal

O exemplo máis coñecido e máis útil de variable aleatoria continua vén dado pola *distribución normal* ou campá de Gauss de media μ e desviación típica σ . Denótase por $N(\mu, \sigma)$ e está definida mediante a función de densidade

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}.$$

A función de densidade da distribución normal está definida e é positiva en toda a recta real. Ademais, é simétrica respecto da súa media.

A distribución normal apareceu como un xeito de estimar as desviacións debidas a erros de medida. Tal propiedade está xustificada matematicamente polo teorema central do límite:

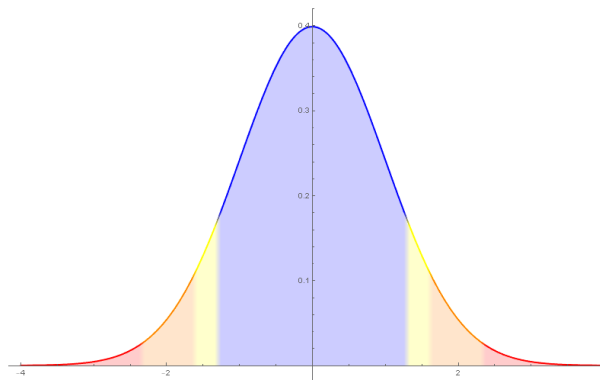


Figura 1.2: Distribuição normal estándar

Teorema 1.4. (Teorema central do limite) *O promedio de moitas variables aleatorias arbitrarias independentes e coa mesma distribución ten, aproximadamente, unha distribución normal.*

Capítulo 2

Introducción á inferencia estadística. Estimación

2.1. Poboación e mostra

A **poboación** é o conxunto de individuos ou obxectos que queremos estudar.

A nosa *hipótese* de partida é que a nosa poboación ten unha característica que pretendemos estudar (por exemplo, estatura, peso, etc.) que segue unha distribución da que coñecemos a súa forma xeral (modelo) pero da que descoñecemos os seus parámetros. Por exemplo, sábese que a estatura segue (aproximadamente) unha distribución normal, pero non coñecemos nin a media nin a desviación típica dunha poboación dada.

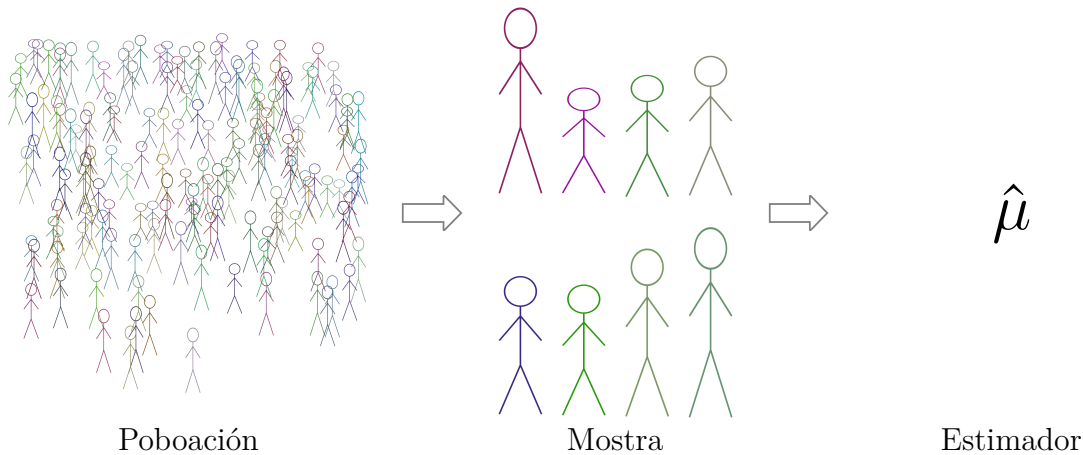
Unha mostra aleatoria é un experimento consistente en tomar n individuos da poboación. Suporemos que a mostra aleatoria se consegue extraendo individuos de xeito *independente*, de modo que tódolos individuos teñan a *mesma probabilidade de ser elexidos en cada momento*. Por tanto, construímos así n variables aleatorias X_1, \dots, X_n independentes e coa mesma distribución de probabilidade cá da poboación. Nótese que despois de face-lo experimento teremos uns valores concretos x_1, \dots, x_n , pero mentres deseñámo-lo experimento eses resultados son descoñecidos e por iso son tratados como *variables aleatorias* en vez de como números; en efecto, antes de realiza-lo experimento estamos extraendo un individuo descoñecido da poboación, e por tanto, a característica que lle estudamos ten a mesma distribución cá da poboación. Dise que n é o **tamaño mostral**, e que X_1, \dots, X_n é unha **mostra aleatoria simple**.

É imposible, sen empregar teoría da probabilidade, decidir de xeito científico o tamaño mostral. Por iso diremos que este é n , e máis adiante intentaremos decidir como se calcula de xeito concreto este valor.

Un **estatístico** é unha función dunha mostra aleatoria simple que expresa unha determinada característica da mostra. Son exemplos de estatísticos a media, a varianza, a cuasivarianza e outras medidas que definimos con anterioridade.

Un **estimador puntual** é un estatístico que toma valores no espazo de parámetros. A súa misión será a de aproximar un parámetro. Un estatístico que ten como misión estimar un parámetro θ denótase $\hat{\theta}$. Por exemplo, se a poboación segue unha distribución normal $N(\mu, \sigma)$, $\hat{\mu}$ será un estimador puntual da media, e $\hat{\sigma}$ un estimador puntual da desviación típica.

Existen varios xeitos de escoller estimadores puntuais. Neste curso non enfatizarémo-la



súa construción, pero si que prestaremos atención a estimadores *insesgados* (aqueles para os que a súa media coincide co valor do parámetro que se pretende estimar) e *consistentes* (aqueles para os que o erro de medida se aproxima a cero cando o tamaño da mostra tende a infinito).

Cando temos uns datos para unha mostra concreta, un estimador puntual dános unha aproximación do parámetro que pretendemos estimar. O problema dun estimador puntual é que non temos idea de se o valor obtido está preto ou lonxe do valor real. Sería interesante ter unha idea do erro cometido coa estimación e acotar probabilisticamente ese erro. Para iso empréganse os chamados intervalos de confianza.

Chámase **intervalo de confianza** a un par de estatísticos T_1 e T_2 , entre os cales se estima que estará certo parámetro descoñecido θ dunha distribución, cunha certa probabilidade de acerto determinada pola condición

$$P\left(T_1(X_1, \dots, X_n) \leq \theta \leq T_2(X_1, \dots, X_n)\right) \geq 1 - \alpha,$$

ou ben,

$$P\left(\theta \in [T_1(X_1, \dots, X_n), T_2(X_1, \dots, X_n)]\right) \geq 1 - \alpha,$$

onde X_1, \dots, X_n é unha mostra aleatoria simple. A probabilidade de éxito na estimación $1 - \alpha$ denomínase **nivel de confianza**. Nestas circunstancias, α é o erro aleatorio ou **nivel de significación**.

Na descripción dun intervalo de confianza fálase de que a probabilidade de que un parámetro estea entre dous estatísticos sexa $1 - \alpha$. Esta é a formulación correcta do problema e o xeito de construí-lo intervalo a nivel teórico. Para datos concretos dunha mostra, os estatísticos transfórmanse en dous valores entre os que se cre que o parámetro buscado está con *confianza* $1 - \alpha$. Insistimos en que para valores concretos se fala de confianza, non de probabilidade. Se por exemplo $\alpha = 0,1$, temos unha confianza do 90 % de que o valor real se atope no intervalo calculado, é dicir, que en 90 de cada 100 mostras o intervalo conterá o valor real. Non se pode falar de probabilidade con datos concretos, xa que non hai variables aleatorias e tódolos valores son xa coñecidos.

2.2. Estimación da media poboacional

O problema que tratamos de resolver nesta sección é o de estima-la media dunha poboación que sabemos que segue unha distribución normal de media μ e desviación típica σ (que en principio son o que queremos estimar). Para iso extraemos unha mostra aleatoria simple X_1, \dots, X_n .

2.2.1. Estimación puntual

Un xeito obvio de estima-la media da poboación é emprega-la media da mostraxe.

Definición 2.1. A **media da mostraxe** é o estimador puntual $\hat{\mu} = \bar{X}$, onde

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i.$$

Como X_1, \dots, X_n teñen a mesma distribución $N(\mu, \sigma)$ e son independentes, temos

$$E(\bar{X}) = \mu, \quad V(\bar{X}) = \frac{\sigma^2}{n}.$$

Debido a estas dúas propiedades, a media mostral é un estimador *insesgado* e *consistente*.

2.2.2. Estimación por intervalos: coñecida a varianza poboacional

Supoñamos que a distribución poboacional segue unha distribución normal $N(\mu, \sigma)$ onde a varianza σ^2 é *coñecida*.

Definición 2.2. Se X_1, \dots, X_n é unha mostra aleatoria simple, entón tomámo-lo estatístico

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim Z$$

que segue unha distribución normal estándar $Z = N(0, 1)$.

Fixemos agora un nivel de significación α (ou un nivel de confianza $1 - \alpha$).

Como a distribución normal é simétrica respecto da media, o noso intervalo de confianza tomarémolo da forma $[\bar{X} - \epsilon, \bar{X} + \epsilon]$, onde ϵ é o *erro* arredor da media que permitimos cometer. Así pois necesitamos

$$P(\mu \in [\bar{X} - \epsilon, \bar{X} + \epsilon]) = 1 - \alpha.$$

Tomámo-lo valor $Z_{\alpha/2}$ para o que $P(Z \geq Z_{\alpha/2}) = \alpha/2$.

Así pois témo-la cadea de igualdades

$$\begin{aligned} 1 - \alpha &= P(|\bar{X} - \mu| \leq \epsilon) \\ &= P\left(-\frac{\epsilon}{\sigma/\sqrt{n}} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq \frac{\epsilon}{\sigma/\sqrt{n}}\right) \\ &= 1 - 2P\left(Z > \frac{\epsilon}{\sigma/\sqrt{n}}\right), \end{aligned}$$

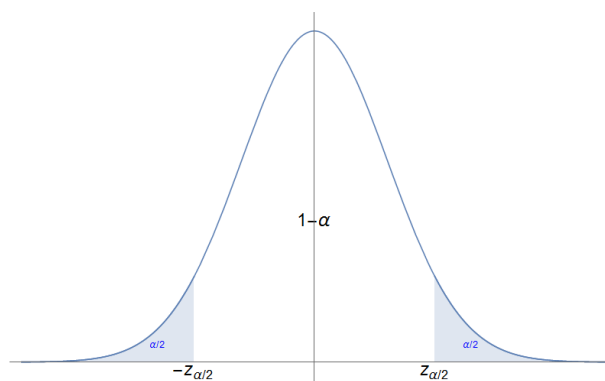


Figura 2.1: Valor para determina-lo intervalo de confianza

de onde se deduce $Z_{\alpha/2} = \frac{\epsilon}{\sigma/\sqrt{n}}$. Despexando ϵ , témo-lo intervalo de confianza

$$\left[\bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right].$$

Observación 2.3. Equivalentemente, resulta máis sinxelo recordar que a partir do estatístico o intervalo de confianza se obtén despexando μ da inecuación

$$\left| \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right| \leq Z_{\alpha/2},$$

ou ben,

$$-Z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq Z_{\alpha/2}.$$

Outro xeito de escribi-lo intervalo de confianza anterior (aproveitando a simetría do mesmo) é mediante a expresión

$$\bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

Problema 2.4. Desexamos estima-lo número medio de latexos por minuto para unha certa poboación. Para iso elíxense aleatoriamente 15 individuos e obtéñense os seguintes resultados:

78	95	70	97	81
85	102	75	78	85
115	80	98	101	92

Supoñendo que a distribución da poboación é normal con desviación típica de 10 latexos por minuto, calcula-lo intervalo de confianza do 99% para a media poboacional de número de latexos por minuto.

Reglas para manipular inecuacións

Sexan x, y números.

Supoñamos $x \leq y$.

Para calquera a ,

$x + a \leq y + a$.

Se $a > 0$, entón $ax \leq ay$.

Se $a < 0$, entón $ax \geq ay$.

Solución. Considerámo-la variable aleatoria X ="número de latexos por minuto". Temos que X ten distribución $N(\mu, 10)$, con μ descoñecido.

En primeiro lugar organizámo-los cálculos para calcula-la media mostral.

X
78
95
70
97
81
85
102
75
78
85
115
80
98
101
92
Σ 1332

Tamaño mostral $n = 15$. Estimación puntual da media $\bar{X} = 1332/15 = 88,8$ latexos.

Nivel de significación: $\alpha = 0,01$. Buscámo-lo valor $Z_{0,005}$ tal que $P(Z \geq Z_{0,005}) = 0,005$. Aproximadamente, $Z_{0,005} = 2,576$.

O intervalo de confianza buscado é entón

$$88,8 \pm 2,576 \cdot \frac{10}{\sqrt{15}} = 88,8 \pm 6,65,$$

que resulta ser $[82,1, 95,5]$.

Conclusión: cunha confianza do 99%, o número medio de latexos por minuto da poboación estudada atópase entre 82.1 e 95.5. \square

Observación 2.5. En ocasións queremos limita-lo erro de estimación para que non sobre-pase certo límite. En tal caso hai que tomar unha mostra suficientemente grande. Como o erro vén dado por $Z_{\alpha/2} \sigma / \sqrt{n}$, se queremos que sexa menor ca ϵ , entón, despexando, obtemos

$$n \geq \left(\frac{Z_{\alpha/2} \sigma}{\epsilon} \right)^2.$$

Observación 2.6. En caso de que a distribución da poboación non se poida garantir que sexa normal, se o tamaño da mostra é grande, o teorema central do límite (Teorema 1.4) dinos que podemos supoñe-la normalidade de \bar{X} , e por tanto, os métodos desta sección seguen sendo aproximadamente válidos. Nos apartados seguintes, se a distribución poboacional non é normal, non se aplica o teorema central do límite aínda que o tamaño da mostra sexa grande, así que neses casos habería que empregar outras técnicas que están máis aló dos obxectivos deste curso.

2.2.3. Estimación por intervalos: descoñecida a varianza poboacional

Supoñamos agora que a distribución poboacional segue unha distribución normal $N(\mu, \sigma)$ onde a varianza σ^2 é descoñecida (o cal é o habitual). Sexa X_1, \dots, X_n é unha mostra aleatoria simple.

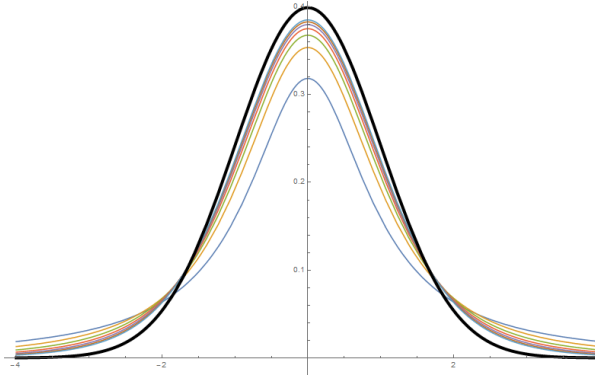


Figura 2.2: Funcións de densidade da t -Student comparadas coa normal estándar

Recordemos que a *cuasi-varianza* ou *varianza mostral* (en contraposición a “varianza poboacional”) vén definida mediante

$$\begin{aligned} s_{n-1}^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \\ &= \frac{1}{n-1} \sum_{i=1}^n X_i^2 - \frac{n}{n-1} \bar{X}^2. \end{aligned}$$

Así, a *cuasi-desviación típica* ou *desviación típica mostral*, s_{n-1} , é a raíz cadrada da cuasi-varianza. Neste curso s denotará, salvo que se diga o contrario, a cuasi-desviación típica s_{n-1} .

Definición 2.7. Para estima-la media cando a varianza poboacional non é coñecida tómasse o estatístico

$$\frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}} \sim t_{n-1}.$$

Este estatístico resulta seguir unha distribución **t -Student de $n-1$ graos de liberdade.**

Definición 2.8. A distribución t -Student é unha nova distribución que ten como función de densidade

$$f(x) = c_{n-1} \left(1 + \frac{t^2}{n-1} \right)^{-\frac{n}{2}},$$

sendo c_{n-1} unha constante que non especificaremos.

Nótese que esta distribución depende dun parámetro n , chamado *graos de liberdade* da distribución, e que haberá que ter en consideración cando mirémo-los valores nas táboas.

Proposición 2.9. *Algunhas propiedades da t -Student:*

- $E(t_n) = 0$ e $V(t_n) = n/(n - 2)$.
- É simétrica respecto da media.
- Ten unha forma parecida á da normal, pero ten cuantiles máis grandes (por tanto produce intervalos de confianza máis grandes).
- Se $n \geq 100$, t_n pode aproximarse por unha $N(0, 1)$.

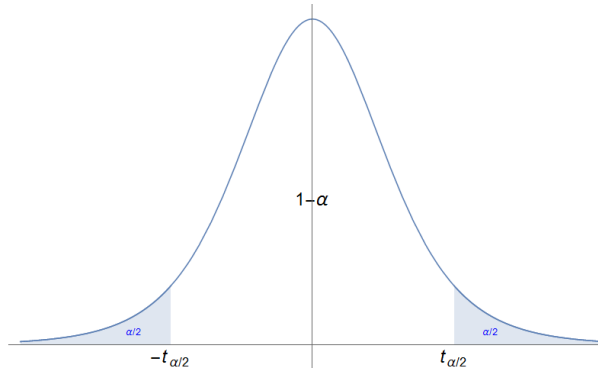


Figura 2.3: Valor para determina-lo intervalo de confianza

Para o cálculo dun intervalo de confianza, o razoamento sería similar ó do anterior apartado. Para un nivel de significación α , o intervalo de confianza para a media vén determinado pola fórmula

$$\left[\bar{X} - t_{n-1, \alpha/2} \frac{s_{n-1}}{\sqrt{n}}, \bar{X} + t_{n-1, \alpha/2} \frac{s_{n-1}}{\sqrt{n}} \right],$$

ou ben,

$$\bar{X} \pm t_{n-1, \alpha/2} \frac{s_{n-1}}{\sqrt{n}},$$

sendo $t_{n-1, \alpha/2}$ o valor tal que $P(t_{n-1} \geq t_{n-1, \alpha/2}) = \alpha/2$.

Observación 2.10. Igual ca no caso anterior, recordar estas fórmulas non resulta sinxelo. Non obstante, coñecido o estatístico necesario para resolve-lo problema, só temos que lembrar que hai que considera-la inecuación

$$-t_{n-1, \alpha/2} \leq \frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}} \leq t_{n-1, \alpha/2},$$

e despegar o valor de μ .

Problema 2.11. Os pesos ó nacer (en gramos) de 10 nenos, elexidos aleatoriamente nun hospital, son:

2750	3316	3969	2211	2806
4195	3061	3827	3572	3430

Supoñendo que a poboación segue unha distribución normal, calcular un intervalo de confianza do 95 % para a media do peso ó nacer dos nenos dese hospital.

Solución. Considerámo-la variable aleatoria X ="peso ó nacer". Temos que X ten distribución $N(\mu, \sigma)$, con μ e σ descoñecidos.

En primeiro lugar, organizámo-los cálculos para a media e cuasi-varianza mostrais.

	X	X^2
	2750	7562500
	3316	10995856
	3969	15752961
	2211	4888521
	2806	7873636
	4195	17598025
	3061	9369721
	3827	14645929
	3572	12759184
	3430	11764900
Σ	33137	113211233

Tamaño mostral $n = 10$. Estimación puntual da media $\bar{X} = 33137/10 = 3313,7$. A cuasi-varianza calcúlase como

$$s_n^2 = \frac{113211233}{10} - 3313,7^2 = 340516,$$

$$s_{n-1}^2 = \frac{10}{9} 340516 = 378351.$$

Extraendo a raíz cadrada obtemos $s_{n-1} = 615,10$.

Nivel de significación: $\alpha = 0,05$. Buscámo-lo valor $t_{9,0,025}$ tal que $P(t_9 > t_{9,0,025}) = 0,025$. Aproximadamente, $t_{9,0,025} = 2,262$.

O intervalo de confianza buscado é entón

$$3313,7 \pm 2,262 \cdot \frac{615,10}{\sqrt{10}} = 3313,7 \pm 440,02,$$

ou explicitamente, $[2873,68, 3753,72]$.

Conclusión: cunha confianza do 95 %, o peso medio ó nacer dos nenos do hospital estudado atópase entre 2873.68 e 3753.72 gramos. \square

2.3. Estimación da varianza poboacional

Nesta sección o problema será o de estima-la varianza dunha poboación que segue unha distribución normal. Tomamos unha mostra aleatoria simple X_1, \dots, X_n .

2.3.1. Estimación puntual

Se a media da poboación é coñecida, tomámo-lo seguinte estimador puntual

Definición 2.12. Definimos

$$s_{\mu}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2.$$

Entón, tense

$$E(s_{\mu}^2) = \sigma^2,$$

é dicir, que s_{μ}^2 é *insesgado*.

Se a media da poboación é descoñecida, o cal é o que sucede habitualmente, cabería pensar que un estimador para a varianza podería ser $s_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$. Isto resulta non se-la mellor idea pois

$$E(s_n^2) = \frac{n-1}{n} \sigma^2,$$

é dicir, que este estimador non é insesgado (ten tendencia a infraestima-la varianza.)

Un xeito máis correcto de estima-la varianza da poboación é emprega-la cuasi-varianza.

Definición 2.13. A *cuasi-varianza* da mostraxe ou *varianza mostral* defínese como

$$s_{n-1}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

Neste caso,

$$E(s_{n-1}^2) = \sigma^2.$$

A cuasi-varianza mostral é un estimador *insesgado*.

2.3.2. Estimación por intervalos: coñecida a media poboacional

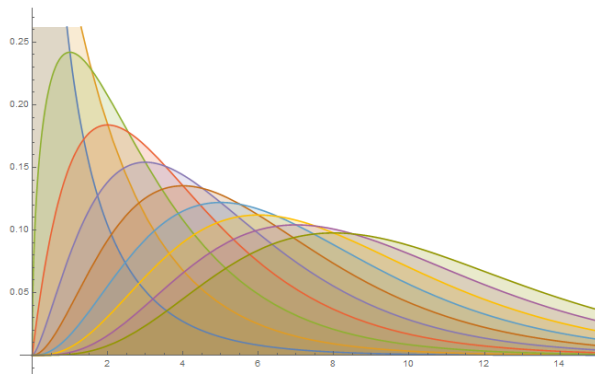


Figura 2.4: Funcións de densidade da χ^2 de Pearson para distintos graos de liberdade

Supoñemos, aínda que normalmente non sucede, que a media poboacional μ é coñecida. É preferible, por tanto, emprega-lo estimador s_{μ} en lugar da cuasi-varianza mostral, xa que o parámetro μ é coñecido exactamente e non cómpre ser aproximado. En realidade esta é unha situación teórica, pois a media poboacional non é coñecida na práctica, pero serve para ir introducindo unha nova distribución que empregaremos máis adiante.

Definición 2.14. Para estima-la varianza poboacional dunha poboación normal con media coñecida toma-lo estatístico

$$\frac{ns_{\mu}^2}{\sigma^2} \sim \chi_n^2,$$

que segue unha distribución χ -cadrado de Pearson con n graos de liberdade.

Definición 2.15. A distribución χ^2 de Pearson ten como función de densidade de probabilidade

$$f(x) = c_n x^{n/2-1} e^{-x/2}, \quad x > 0,$$

onde c_n é unha constante.

A distribución χ^2 de Pearson, ó igual que sucedía coa distribución t de Student, depende dun parámetro que se coñece como o número de graos de liberdade da distribución.

Proposición 2.16. *Algunhas propiedades da χ^2 de Pearson:*

- $E(\chi_n^2) = n$ e $V(\chi_n^2) = 2n$.
- Só está definida para valores positivos e non é simétrica.
- Se $n > 30$, χ_n^2 pode aproximarse por unha normal $N(n, \sqrt{2n})$; unha aproximación aínda mellor é $\sqrt{2\chi_n^2} - \sqrt{2n-1} \cong N(0, 1)$.

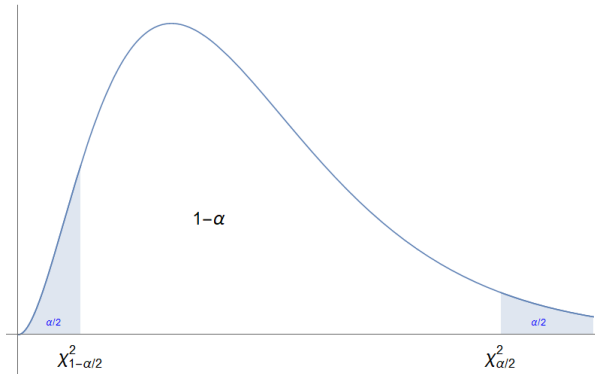


Figura 2.5: Valores para determina-lo intervalo de confianza

Dado que a distribución χ^2 de Pearson non é simétrica, o intervalo de confianza que construímos tampouco o será. Fixado un nivel de significación α , buscamos dous extremos de intervalo a e b de xeito que á esquerda de a e á dereita de b quede probabilidade $\alpha/2$. É dicir, buscámo-los valores $a = \chi_{n,1-\alpha/2}^2$ e $b = \chi_{n,\alpha/2}^2$ tales que $P(\chi_n^2 \geq \chi_{n,1-\alpha/2}^2) = 1 - \alpha/2$ e $P(\chi_n^2 \geq \chi_{n,\alpha/2}^2) = \alpha/2$.

Nestas condicións, o intervalo de confianza para a varianza poboacional buscado vén dado pola fórmula

$$\left[\frac{ns_{\mu}^2}{\chi_{n,\alpha/2}^2}, \frac{ns_{\mu}^2}{\chi_{n,1-\alpha/2}^2} \right].$$

Observación 2.17. Como sempre, resulta máis sinxelo, coñecido o estatístico necesario para estima-la varianza poboacional, calcula-lo intervalo de confianza a partir de despear σ^2 da inecuación

$$\chi_{n,1-\alpha/2}^2 \leq \frac{ns_{\mu}^2}{\sigma^2} \leq \chi_{n,\alpha/2}^2,$$

2.3.3. Estimación por intervalos: descoñecida a media poboacional

O procedemento é similar ó caso anterior, pero agora temos que emprega-la cuasi-varianza mostral.

Definición 2.18. Para estima-la varianza poboacional dunha poboación normal con media descoñecida tomámo-lo estatístico

$$\frac{(n-1)s_{n-1}^2}{\sigma^2} \sim \chi_{n-1}^2,$$

que segue unha distribución χ^2 de Pearson con $n-1$ graos de liberdade.

Observación 2.19. O procedemento para atopar un intervalo de confianza é similar a casos anteriores. De feito, o intervalo de confianza buscado, para unha nivel de significación α , é determinado por

$$\chi_{n-1,1-\alpha/2}^2 \leq \frac{(n-1)s_{n-1}^2}{\sigma^2} \leq \chi_{n-1,\alpha/2}^2$$

Despejando σ^2 obtemos:

$$\left[\frac{(n-1)s_{n-1}^2}{\chi_{n-1,\alpha/2}^2}, \frac{(n-1)s_{n-1}^2}{\chi_{n-1,1-\alpha/2}^2} \right].$$

Problema 2.20. Obtense unha mostra aleatoria de 100 adultos aparentemente sans co fin de establecer un patrón con respecto ó que se considerará unha lectura normal de calcio. Extráese unha mostra de sangue de cada adulto. A variable estudada é X ="contido de calcio en mg/dl de sangue", que se supón que presenta unha distribución aproximadamente normal. Obtívose unha media mostral de 9.5mg/dl e unha varianza $s_n^2 = 0,2475$. Calcular intervalos de confianza do 99% para a media e a desviación típica da poboación.

Solución. Considerámo-la variable aleatoria X ="contido de calcio en mg/dl de sangue".

Os datos que temos no enunciado son o tamaño da mostra $n = 100$, a media mostral $\bar{X} = 9,5$ e a varianza $s_n^2 = 0,2475$. A cuasi-varianza é $s_{n-1}^2 = \frac{100}{99} \cdot 0,2475 = 0,25$; logo $s_{n-1} = 0,5$. O nivel de significación é $\alpha = 0,01$.

Para o cálculo dun intervalo de confianza para a media buscámo-lo valor $t_{99,0,005} = 2,63$. Así un intervalo para a media é

$$9,5 \pm 2,63 \cdot \frac{0,5}{\sqrt{100}} = 9,5 \pm 0,13,$$

ou ben, $[9,37, 9,63]$.

A continuación pasamos á varianza. Temos que buscar *dous* valores da χ^2 : $\chi_{99,0,005}^2 = 138,99$ e $\chi_{99,0,995}^2 = 66,51$. O intervalo de confianza para a varianza é

$$\left[\frac{99 \cdot 0,25}{138,99}, \frac{99 \cdot 0,25}{66,51} \right] = [0,18, 0,37].$$

Simplemente extraendo raíces cadradas temos un intervalo de confianza para a desviación típica: $[0,42, 0,61]$.

Conclusión: cunha confianza do 99 %, o contido en calcio en sangue medido en mg/dl na poboación estudada ten unha media que está comprendida entre 9.37 e 9.63, e unha desviación típica entre 0.42 e 0.61. \square

2.4. Estimación dunha proporción

Supoñamos que temos unha variable con dous posibles valores. Temos unha poboación na que queremos estima-la proporción p de individuos que teñen un deses valores. Unha mostra individual desa poboación seguirá pois unha distribución de Bernoulli de parámetro p , mentres que a poboación segue unha distribución *binomial* de parámetros N (número de elementos) e p .

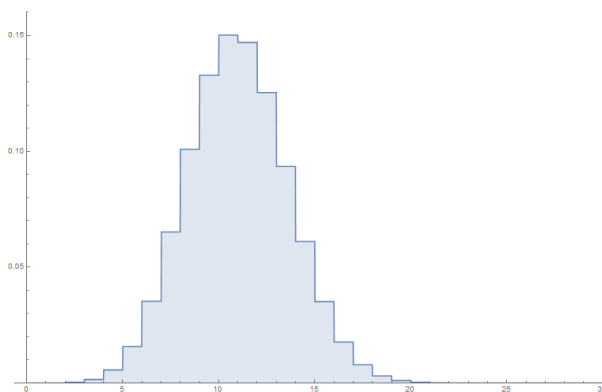


Figura 2.6: Función de masas dunha binomial (30, 0.35)

Recordemos que a distribución binomial de parámetros N e p é unha distribución discreta con función de masa

$$P(X = k) = \binom{N}{k} p^k (1 - p)^{N-k}.$$

A súa media e a súa varianza son

$$E(X) = Np, \quad V(X) = Np(1 - p).$$

2.4.1. Estimación puntual

Queremos construír un estimador \hat{p} de p . Para iso definímo-la variable aleatoria X que lle asigna 1 ó valor que queremos medir, e 0 ó outro. Escollemos unha mostra aleatoria simple X_1, \dots, X_n .

Definición 2.21. Para estimar unha proporción é razoable toma-lo estimador puntual

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i,$$

que aproxima a proporción da característica que queremos medir cos datos da mostra escollida.

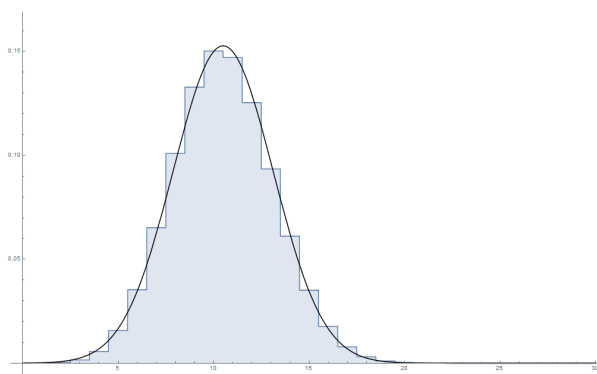


Figura 2.7: A anterior distribución binomial comparada cunha normal da mesma media e varianza

Temos que $n\hat{p} = \sum_{i=1}^n X_i$ segue unha distribución binomial de parámetros n e p . No caso de que a mostra sexa grande (con $np, n(1-p) \geq 5$ acostuma ser suficiente), podemos aproxima-la binomial por unha normal.

Co obxectivo de estandariza-los cálculos e facer máis inmediato o emprego das táboas realizaremos o procedemento típico de tipifica-la variable.

Definición 2.22. Por tanto, habitualmente consideraremos que a distribución na mostraxe para estimar unha proporción vén dada por

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim Z$$

que segue (aproximadamente) unha $N(0, 1)$.

Satisfaise que

$$E(\hat{p}) = p, \quad V(\hat{p}) = \frac{p(1-p)}{n},$$

e por tanto, dise que \hat{p} é un estimador *insesgado* e *consistente* de p .

2.4.2. Estimación por intervalos

O procedemento para atopar un intervalo de confianza é similar ó explicado para a media, aínda que hai algunha dificultade que presentamos a continuación. Sexa α o nivel de significación. Tomamos $Z_{\alpha/2}$ tal que $P(z \geq Z_{\alpha/2}) = \alpha/2$. En principio o cálculo dun intervalo de confianza viría expresado despegando p na fórmula

$$\left| \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \right| \leq Z_{\alpha/2}.$$

O problema é que o denominador $\sqrt{p(1-p)/n}$ depende de p , que é xusto o que queremos estimar. En consecuencia, aproximaremos $\sqrt{p(1-p)/n}$ por $\sqrt{\hat{p}(1-\hat{p})/n}$.

Observación 2.23. Así un intervalo de confianza para a proporción vén dado pola expresión

$$-Z_{\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq Z_{\alpha/2}.$$

Despexando p obtemos.

$$\left[\hat{p} - Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right],$$

ou ben,

$$\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}.$$

Problema 2.24. Un laboratorio desexa averiguar a proporción de cápsulas defectuosas que produce dun determinado medicamento. Para iso selecciona e proba 2000 unidades e descubre un total de 200 unidades defectuosas. Estima a proporción de cápsulas defectuosas na produción. Calcular un intervalo de confianza ó 95 % para a proporción.

Solución. Considerámo-la variable aleatoria X que asigna o valor 1 ás cápsulas defectuosas e 0 ás correctas.

Tamaño mostral $n = 2000$. Estimación puntual da proporción $\hat{p} = 200/2000 = 0,1 = 10\%$.

Nivel de significación: $\alpha = 0,05$. Buscámo-lo valor $Z_{0,025}$ tal que $P(Z \geq Z_{0,025}) = 0,025$. Aproximadamente, $Z_{0,025} = 1,96$.

O intervalo de confianza buscado é entón

$$0,1 \pm 1,96 \sqrt{\frac{0,1(1-0,1)}{2000}} = 0,1 \pm 0,0131,$$

que explicitamente, en termos de porcentaxes, é [8,69 %, 11,31 %].

Conclusión: cunha confianza do 95 %, a porcentaxe de cápsulas defectuosas na produción do laboratorio sitúase entre o 8,96 % e o 11,31 %. \square

Observación 2.25. Determinación do tamaño da mostra

En vista do intervalo de confianza construído para a proporción, o erro cometido ó tomar \hat{p} en lugar do valor verdadeiro p estímase que é

$$Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}},$$

que depende do tamaño mostral n , do nivel de confianza α , e de $\sqrt{\hat{p}(1-\hat{p})}$. Se coñecemos (ou podemos estimar con precisión) o valor de \hat{p} , bastaría impoñer que a anterior fórmula é $< \epsilon$ e despexar n .

Cando o valor de \hat{p} non é coñecido pode estimarse o tamaño da mostra necesario para limita-lo erro, se ben o valor obtido será máis grande que cando \hat{p} é coñecido. No intervalo $[0, 1]$ pode verse, empregando as técnicas do cálculo, que o máximo de $\sqrt{x(1-x)}$ está en $x = 1/2$, de xeito que teremos sempre

$$Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq Z_{\alpha/2} \sqrt{\frac{0,5(1-0,5)}{n}}.$$

Se queremos que o erro sexa menor ca ϵ , basta entón impoñe-la condición

$$Z_{\alpha/2} \sqrt{\frac{0,5(1-0,5)}{n}} < \epsilon,$$

de onde resulta

$$n > \frac{Z_{\alpha/2}^2}{4\epsilon^2}.$$

Problema 2.26. Para toma-la decisión de someter ou non a referendo unha lei, o goberno dun certo país necesita encargar un estudo sobre a porcentaxe de votantes que a apoiaría. Dada a importancia política da mesma e a polémica xurdida, necesita unha estimación do voto cun erro menor do 1%. ¿Cal sería o tamaño mostral mínimo requerido para un nivel de confianza do 99%?

Solución. Considerámo-la variable aleatoria X ="intención de voto".

Para estima-lo tamaño da mostra para unha proporción (Observación 2.25), empregámo-lo estatístico

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}},$$

que segue unha distribución normal estándar. Despexando p da desigualdade

$$\left| \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \right| \leq Z_{\alpha/2},$$

obtense a fórmula

$$\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}.$$

A estimación do erro é

$$Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}.$$

Neste caso non temos unha estimación da proporción \hat{p} . É sinxelo ver que a función $x \mapsto \sqrt{x(1-x)}$ alcanza o seu máximo no intervalo $[0, 1]$ no punto $x = 1/2$. Por tanto, necesitamos despexar n da desigualdade $Z_{\alpha/2} \sqrt{\frac{0,5(1-0,5)}{n}} \leq \epsilon$, onde ϵ é o valor fixado polo problema. Así, obtense $n \geq \left(\frac{Z_{\alpha/2}}{2\epsilon}\right)^2$.

O nivel de significación é $\alpha = 0,01$. Calculamos $Z_{0,005} = 2,5758$. Neste caso $\epsilon = 0,01$. Substituíndo na fórmula, $n \geq \left(\frac{2,5758}{2 \cdot 0,01}\right)^2 = 16587,2415$.

Conclusión: para que a diferenza entre a proporción mostral e a proporción poboacional de intención de voto sexa como moito de $\pm 0,01\%$ cun nivel de confianza do 99,0%, teríamos que tomar unha mostra de polo menos 16588 persoas. \square

2.5. Resumo de estimadores

Táboa resumo cos resultados explicados neste capítulo.

Estimación da media poboacional	
Varianza poboacional coñecida	
Estimador puntual	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
Distribución na mostraxe	$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim Z = N(0, 1)$
Inecuación	$-Z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq Z_{\alpha/2}$
Intervalo de confianza	$\left[\bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$
Valores en táboa	$P(N(0, 1) \geq Z_{\alpha/2}) = \frac{\alpha}{2}$
Varianza poboacional descoñecida	
Estimador puntual	$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$
Distribución na mostraxe	$\frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}} \sim t_{n-1}$
Inecuación	$-t_{n-1, \alpha/2} \leq \frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}} \leq t_{n-1, \alpha/2}$
Intervalo de confianza	$\left[\bar{X} - t_{n-1, \alpha/2} \frac{s_{n-1}}{\sqrt{n}}, \bar{X} + t_{n-1, \alpha/2} \frac{s_{n-1}}{\sqrt{n}} \right]$
Valores en táboa	$P(t_{n-1} \geq t_{n-1, \alpha/2}) = \frac{\alpha}{2}$

Estimación da varianza poboacional

Media poboacional coñecida

Estimador puntual	$s_{\mu}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$
Distribución na mostraxe	$\frac{ns_{\mu}^2}{\sigma^2} \sim \chi_n^2$
Inecuación	$\chi_{n,1-\alpha/2}^2 \leq \frac{ns_{\mu}^2}{\sigma^2} \leq \chi_{n,\alpha/2}^2$
Intervalo de confianza	$\left[\frac{ns_{\mu}^2}{\chi_{n,\alpha/2}^2}, \frac{ns_{\mu}^2}{\chi_{n,1-\alpha/2}^2} \right]$
Valores en táboa	$P(\chi_n^2 \geq \chi_{n,\alpha/2}^2) = \frac{\alpha}{2}$ $P(\chi_n^2 \geq \chi_{n,1-\alpha/2}^2) = 1 - \frac{\alpha}{2}$

Media poboacional descoñecida

Estimador puntual	$s_{n-1}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$
Distribución na mostraxe	$\frac{(n-1)s_{n-1}^2}{\sigma^2} \sim \chi_{n-1}^2$
Inecuación	$\chi_{n-1,1-\alpha/2}^2 \leq \frac{(n-1)s_{n-1}^2}{\sigma^2} \leq \chi_{n-1,\alpha/2}^2$
Intervalo de confianza	$\left[\frac{(n-1)s_{n-1}^2}{\chi_{n-1,\alpha/2}^2}, \frac{(n-1)s_{n-1}^2}{\chi_{n-1,1-\alpha/2}^2} \right]$
Valores en táboa	$P(\chi_{n-1}^2 \geq \chi_{n-1,\alpha/2}^2) = \frac{\alpha}{2}$ $P(\chi_{n-1}^2 \geq \chi_{n-1,1-\alpha/2}^2) = 1 - \frac{\alpha}{2}$

Estimación dunha proporción

Estimador puntual	$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i$
Distribución na mostraxe	$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim Z = N(0, 1)$
Inecuación	$-Z_{\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq Z_{\alpha/2}$
Intervalo de confianza	$\left[\hat{p} - Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$
Valores en táboa	$P(N(0, 1) \geq Z_{\alpha/2}) = \frac{\alpha}{2}$

Capítulo 3

Contraste de hipóteses

A finalidade do contraste de hipóteses é decidir se unha determinada hipótese ou afirmación sobre a distribución da poboación pode ser invalidada estatisticamente a partir das observacións contidas nunha mostra.

A hipótese sobre a distribución da poboación denomínase xenericamente **hipótese nula** e désígnase por H_0 . Esta pretende contrastarse fronte a unha segunda hipótese chamada **hipótese alternativa** H_1 , que agrupa a tódalas posibles poboacións nas que H_0 non é certa.

O contraste de hipóteses non ten normalmente un comportamento imparcial fronte a H_0 e H_1 , xa que o problema consiste, non en decidir cal das dúas suposicións é máis verosímil en vista dos datos, senón en decidir se a mostra proporciona ou non evidencia suficiente para descartar H_0 en favor de H_1 .

Nun problema de contraste de hipóteses os *dous únicos resultados posibles* consisten en *rexeitar* H_0 ou non rexeitar (ou aceptar) H_0 . En xeral, o obxectivo cando se fai un contraste de hipóteses é tratar de rexeitar H_0 , é dicir, de intentar dar evidencia estatística suficiente para concluír que a hipótese alternativa H_1 é certa. Por exemplo, se queremos probar estatisticamente que un determinado medicamento é útil para curar unha enfermidade, a nosa hipótese nula H_0 será formular matematicamente que o medicamento non é útil, e a nosa hipótese alternativa, que si que o é. Se conseguimos rexeitar H_0 teremos probado estatisticamente que o medicamento é útil. En caso contrario, aceptaremos H_0 e concluiremos que non hai evidencia de que o medicamento en cuestión sirva para cura-la enfermidade.

A decisión de rexeitar ou non H_0 deberá facerse en vista dos valores obtidos nunha mostra dalgún estatístico que ten unha distribución de probabilidade que, baixo a presunción de que H_0 é certa, é coñecido. Este estatístico denomínase **estatístico de contraste**.

Por tanto, un contraste de hipóteses consiste en dividi-lo espazo mostral en dúas rexións disxuntas. Unha dela chámase **rexión crítica** ou de rexeitamento, e se a mostra pertence a ela, rexéitase H_0 para inclinarse por H_1 . A outra chámase **rexión de aceptación**, na que H_0 é aceptada en caso de que a mostra pertenza a ela.

Tal e como está presentado o problema existen dúas disxuntivas: a veracidade ou falsidade da hipótese nula, e aceptar ou rexeitar esta. Así, temos a seguinte casuística:

	H_0 é certa	H_0 é falsa
rexeitar H_0	erro de tipo I	decisión correcta
aceptar H_0	decisión correcta	erro de tipo II

Observación 3.1. Para o contraste de hipóteses resulta ás veces interesante facer un similar co sistema xurídico americano. O veredicto dun xurado con respecto a un crime ten dúas posibles decisións: “culpable” ou “non culpable”. Nunca se dictamina que alguén é “inocente”: a inocencia presuponse (hipótese nula H_0) e non é necesario probala. O que si é necesario probar é a culpabilidade (hipótese alternativa H_1). Neste sistema intenta minimizase que os inocentes sexan condenados (erro de tipo I), aínda a costa de que haxa culpables que queden impunes (erro de tipo II).

Como en xeral é imposible minimizar simultaneamente os tipos de erro I e II, o criterio tradicional na teoría de contrastes consiste en:

1. Fixar un límite para a probabilidade de cometer un erro de tipo I, chamado **nivel de significación**

$$\alpha = P(\text{rexeitar } H_0 \mid H_0 \text{ é certa}).$$

A $1 - \alpha$ chámasele nivel de confianza.

2. Rexeitar todos aqueles tests que imponen que a probabilidade de rexeitar H_0 cando sexa certa non supere o valor α do nivel de significación.
3. Entre tódolos test non excluídos anteriormente, tratar de minimiza-la probabilidade de erro de tipo II. Chámase **potencia** á probabilidade de detectar que unha hipótese é falsa,

$$\begin{aligned} \beta &= P(\text{rexeitar } H_0 \mid H_0 \text{ é falsa}) \\ &= 1 - P(\text{erro de tipo II}), \end{aligned}$$

e por tanto preténdese maximiza-la potencia do método.

Tal procedemento outorga, en principio, prioridade a rebaixa-lo risco de erro de tipo I por debaixo do nivel de significación. De aí que o tratamento que reciben ambas hipóteses sexa asimétrico e estas non sexan intercambiáveis. De feito, no contraste de hipóteses considérase que H_0 é a hipótese establecida, que ten presunción de veracidade, e contra a cal é necesario esgrimir unha grande evidencia para poder invalidala. Así, emprégase un carácter conservador a favor da hipótese H_0 : o nivel de significación que se fixa intenta garantir que sexa moi infrecuente rexeita-la hipótese correcta. A preocupación por deixar vixente unha hipótese nula falsa (erro de tipo II) é menor, polo que pode aceptarse nese caso un risco máis alto. En consecuencia, se o resultado dun contraste de hipóteses é acepta-la hipótese nula, debe interpretarse que as observacións non aportaron suficiente evidencia para descartala. Pola contra, se se rexeita é porque se está razoablemente seguro de que H_0 é falsa e H_1 é verdadeira.

O rango de valores α debe estar adaptado á importancia ou trascendencia do problema. A elección do nivel de significación é unha cuestión delicada e importante á que se lle debe prestar atención. Fixémonos cal é a razón de chamarlle a α “nivel de significación”. Cando rexeitamos a hipótese nula, é porque obtivemos unha mostra que dá evidencia clara de que esta é falsa. Aínda cabería a posibilidade de que a mostra elixida fose “mala”, no sentido de que non representa realmente a poboación. Non obstante, a probabilidade de que iso sucedese é menor ca α , e por tanto considérase moi improbable: é difícil que tal mostra aporte eses datos como consecuencia razoable das fluctuacións aleatorias debidas á súa elección. En consecuencia, decídese que a mostra é *significativa*, e rexéitase a hipótese nula.

Nos problemas estatísticos *paramétricos* nos que a distribución da poboación pertence a unha familia con parámetros nun conxunto Θ , tanto a hipótese nula como a alternativa serán especificadas mediante subconxuntos disxuntos Θ_0 e Θ_1 tales que $\Theta_0 \cup \Theta_1 = \Theta$. Deste xeito o contraste de hipóteses escríbese como

$$H_0: \theta \in \Theta_0,$$

$$H_1: \theta \in \Theta_1.$$

Habitualmente os contrastes de hipóteses estudados correspóndense con dúas posibles situacións: os **contrastos bilaterais** que nós tomaremos da forma $H_0: \theta = \theta_0$, $H_1: \theta \neq \theta_0$, e os **contrastos unilaterais** $H_0: \theta \leq \theta_0$, $H_1: \theta > \theta_0$ (ou coas desigualdades invertidas).

Os métodos estatísticos para o deseño de tests de hipóteses son complicados e están fóra dos obxectivos deste curso. Non obstante presentaremos os procedementos para realizar contrastes de hipóteses para poboacións normais nos que se contrastan as características máis habituais.

3.1. Contraste de hipóteses para a media da poboación

Supoñamos que temos unha determinada poboación que se rixe por unha distribución de probabilidade normal. Temos unha certa suposición sobre a media e queremos contrasta-la súa veracidade. Para iso tomamos unha mostra aleatoria simple X_1, \dots, X_n .

3.1.1. Contrastos bilaterais

Empezamos co caso en que contrastamos un determinado valor da media. Así,

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

Supoñendo que H_0 fose certa, tomámo-lo estatístico de contraste

$$\frac{\bar{X} - \mu_0}{s_{n-1}/\sqrt{n}} \sim t_{n-1}$$

que ten, como vimos na sección dedicada ó cálculo de intervalos de confianza para a media con varianza descoñecida (Definición 2.7), unha distribución *t*-Student (Definición 2.8) con $n - 1$ graos de liberdade. (En caso de que a varianza poboacional σ fose coñecida tomaríamo-lo estatístico $\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim Z$, que ten distribución $Z = N(0, 1)$, como consta na sección dedicada ó cálculo de intervalos de confianza para a media con varianza coñecida (Definición 2.2).)

Tomamos un nivel de significación α .

- A rexión crítica é $(-\infty, -t_{n-1, \alpha/2}) \cup (t_{n-1, \alpha/2}, +\infty)$, é dicir, cando

$$\left| \frac{\bar{X} - \mu_0}{s_{n-1}/\sqrt{n}} \right| > t_{n-1, \alpha/2}.$$

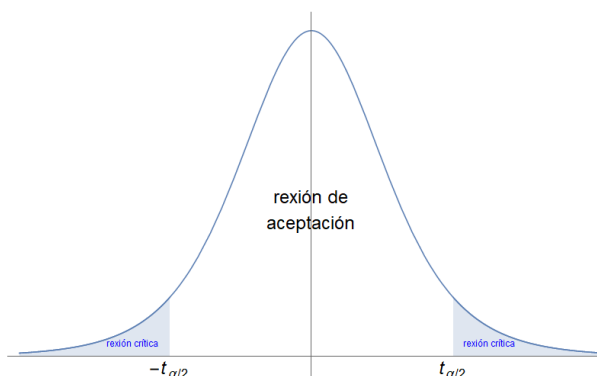


Figura 3.1: Región de aceptación para un contraste bilateral

- A rexión de aceptación é por tanto $[-t_{n-1, \alpha/2}, t_{n-1, \alpha/2}]$.

(En caso de que a varianza sexa coñecida, a rexión crítica é $(-\infty, -Z_{\alpha/2}) \cup (Z_{\alpha/2}, +\infty)$, e a rexión de aceptación é $[-Z_{\alpha/2}, Z_{\alpha/2}]$.)

Cando o valor obtido na mostra está dentro do intervalo de aceptación, aceptamos H_0 . Cando está na rexión crítica, é dicir, fóra do intervalo de aceptación, rexeitamos H_0 . Neste caso sempre existe a pequena posibilidade α de que a mostra tomada non sexa representativa da poboación e cometamos un erro de tipo I (rexeitar un modelo correcto); non obstante, a probabilidade disto é pequena, e en vista dos datos deberemos de rexeita-la hipótese nula.

Problema 3.2. Estudámo-lo crecemento anual dos abetos. Cremos que o valor medio desta variable é $\mu_0 = 7,25$. Non obstante, nunha mostra de 50 árbores obtívose o valor $\bar{X} = 7,27$ e $s_{n-1} = 0,03$. ¿É este resultado compatible coa nosa suposición cun nivel de confianza do 95%?

Solución. Estudámo-la variable aleatoria $X = \text{“crecemento anual dos abetos”}$.

Neste caso témo-lo contraste de hipóteses

$$H_0: \mu = 7,25, \quad H_1: \mu \neq 7,25.$$

O nivel de significación é $\alpha = 0,05$. Damos como datos $n = 50$, $\bar{X} = 7,27$, $s_{n-1} = 0,03$.

Xa que a varianza da poboación non é coñecida, empregamos un estatístico $\frac{\bar{X} - \mu_0}{s_{n-1}/\sqrt{n}} \sim t_{n-1}$ que ten distribución t -Student, e obtemos $t_{49, 0,025} = 2,01$. Como

$$\frac{7,27 - 7,25}{0,03/\sqrt{50}} = 4,71 \notin [-2,01, 2,01],$$

o valor obtido está fóra do intervalo de aceptación.

Conclusión: rexeitamos H_0 e deducimos que hai evidencia significativa, polo menos do 95%, de que o valor medio de crecemento anual dos abetos non é $\mu_0 = 7,25$ metros. \square

3.1.2. Contrastes unilaterais

Neste caso a hipótese nula establece un límite superior ou inferior para a media. Así escribiremos

$$H_0: \mu \leq \mu_0, \quad H_1: \mu > \mu_0,$$

que é un *contraste unilateral dereito*, ou ben,

$$H_0: \mu \geq \mu_0, \quad H_1: \mu < \mu_0,$$

para un *contraste unilateral esquerdo*.

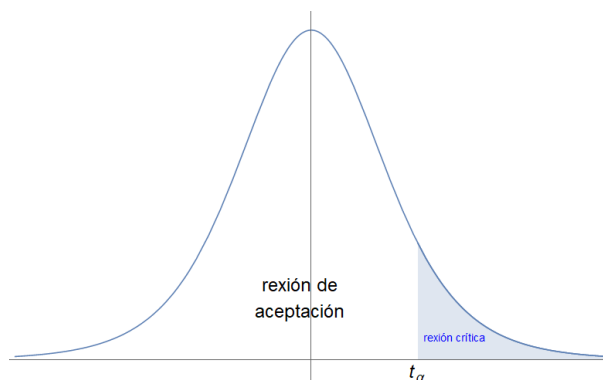


Figura 3.2: Rexión de aceptación para un contraste unilateral dereito

De novo, tomámo-lo estatístico

$$\frac{\bar{X} - \mu_0}{s_{n-1}/\sqrt{n}} \sim t_{n-1},$$

que ten, unha distribución *t*-Student con $n - 1$ graos de liberdade (Definición 2.8).

Tomamos un nivel de significación α .

- A rexión crítica é $(t_{n-1, \alpha}, +\infty)$ para un contraste unilateral dereito, e $(-\infty, -t_{n-1, \alpha})$ para un contraste unilateral esquerdo.
- A rexión de aceptación é por tanto $(-\infty, t_{n-1, \alpha}]$ para un contraste unilateral dereito, e $[-t_{n-1, \alpha}, +\infty)$ para un contraste unilateral esquerdo.

(En caso de que a varianza poboacional σ fose coñecida (Definición 2.2) tomaríámolo estatístico $\frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim Z$, que ten distribución normal $Z = N(0, 1)$. A rexión crítica é $(Z_\alpha, +\infty)$ para o contraste unilateral dereito, e $(-\infty, -Z_\alpha)$ para o contraste unilateral esquerdo. As rexións de aceptación son, respectivamente, $(-\infty, Z_\alpha]$ e $[-Z_\alpha, +\infty)$.)

Igual ca no caso anterior, cando o valor obtido na mostra está dentro do intervalo de aceptación, aceptamos H_0 . Cando está na rexión crítica, é dicir, fóra do intervalo de aceptación, rexeitamos H_0 .

Problema 3.3. A consellería de pesca considera que non se deben extraer ameixas se o número medio de bacterias por centímetro cúbico na auga sobrepasa 70. Como norma xeral, as rías galegas están por debaixo dese nivel de concentración. Fíxose unha mostraxe en 9 lugares da ría e obtívose un reconto $\bar{X} = 71,7$, con $s_{n-1} = 2,3$. ¿Que decisión deben toma-los inspectores con nivel de confianza 99%?

Solución. Considérase a variable aleatoria X ="número medio de bacterias por centímetro cúbico na auga".

O contraste de hipóteses a considerar é

$$H_0: \mu \leq 70, \quad H_1: \mu > 70.$$

O nivel de significación é $\alpha = 0,01$, e temos $\mu_0 = 70$, $n = 9$, $\bar{X} = 71,7$, $s_{n-1} = 2,3$.

Empregámo-lo estatístico $\frac{\bar{X}-\mu_0}{s_{n-1}/\sqrt{n}}$, que ten distribución t -Student, e obtemos $t_{8,0,01} = 2,896$. Como

$$\frac{71,7 - 70}{2,3/\sqrt{9}} = 2,22 \in (-\infty, 2,896],$$

o valor obtido está dentro do intervalo de aceptación. Así, este número anormalmente alto non é significativo e probablemente se deba á elección da mostra.

Conclusión: aceptamos H_0 , o cal quere dicir que non hai evidencia significativa, polo menos do 99%, de que o número medio de bacterias por centímetro cúbico de auga é menor ou igual ca 70. En consecuencia as ameixas son aptas para o consumo. \square

Para a elección dunha hipótese nula nun contraste unilateral debe considerarse aquela desigualdade para a que se desexa minimiza-la probabilidade de erro de tipo I (rexeitar H_0 sendo certa). É dicir, que H_0 é a hipótese contra a que hai que esgrimir unha evidencia contundente para rexeitala. Recordemos que nun contraste de hipóteses aquilo que queremos probar debe estar contido na hipótese nula.

Problema 3.4. A normativa cambia e a consellería de pesca require evidencia significativa de que o número medio de bacterias por centímetro cúbico na auga sexa menor ca 70 para permiti-la extracción de ameixas; é fundamental asegurarse de que tal número non é sobrepasado. Coa mostra de 9 lugares da ría obtida de $\bar{X} = 71,7$, e $s_{n-1} = 2,3$, ¿que decisión deben toma-los inspectores con nivel de confianza 99%? ¿E se fose $\bar{X} = 68,7$?

Solución. A variable aleatoria considerada segue sendo X ="número medio de bacterias por centímetro cúbico na auga".

Como agora é importante non sobrepasa-lo valor 70, e cómpre dar evidencia concluínte diso, o contraste de hipóteses a considerar é

$$H_0: \mu \geq 70, \quad H_1: \mu < 70.$$

Igual ca antes, o nivel de significación é $\alpha = 0,01$, e temos $\mu_0 = 70$, $n = 9$, $\bar{X} = 71,7$, $s_{n-1} = 2,3$.

Empregámo-lo estatístico $\frac{\bar{X}-\mu_0}{s_{n-1}/\sqrt{n}}$, que ten distribución t -Student, e obtemos $t_{8,-0,01} = 2,896$. Como

$$\frac{71,7 - 70}{2,3/\sqrt{9}} = 2,22 \in [-2,896, +\infty),$$

o valor obtido está dentro do intervalo de aceptación.

Conclusión: *aceptamos* H_0 , co que non hai evidencia significativa cunha confianza do 99 % de que o número medio de bacterias por centímetro cúbico da auga sexa menor ou igual ca 70. Por tanto, hai que *prohibi-la extracción de ameixa*.

Nota: en realidade resulta superfluo facer un contraste de hipóteses para este caso, xa que a mostra non dá evidencia en contra da hipótese nula (satisfai $\bar{X} \leq \mu_0 = 70$). Non obstante, vemos que os cálculos claramente confirman esta afirmación.

Para $\bar{X} = 68,7$ teriamos

$$\frac{68,7 - 70}{2,3/\sqrt{9}} = -1,696 \in [-2,896, +\infty),$$

co que aínda neste caso *aceptamos* H_0 .

Conclusión: *aceptamos* H_0 , co que non hai evidencia significativa cunha confianza do 99 % de que o número medio de bacterias por centímetro cúbico da auga sexa menor ou igual ca 70. Por tanto, tamén neste caso habería que *prohibi-la extracción de ameixa*.

Nótese que neste caso é importante que o nivel medio de bacterias sexa menor ca 70, e por tanto é necesario asegurarse que un valor medio pequeno na mostra non é froito do azar ó escollela. \square

3.1.3. O valor P ou valor crítico

Intuitivamente o valor P é un número que dá o grao de sorpresa que un experimento causaría nun partidario da hipótese nula.

Definición 3.5. Para un contraste unilateral dereito correspóndese coa área baixo a curva da función de densidade dunha variable aleatoria X cara á dereita do valor observado polo estatístico de contraste, é dicir,

$$P = P(X \geq \text{valor no estatístico}).$$

Para un contraste unilateral esquerdo o **valor** P é a área baixo a curva da función de densidade cara á esquerda do valor observado polo estatístico de contraste.

Por tanto, cando fagamos un contraste de hipóteses rexeitaremos H_0 cando creamos que o valor P é demasiado pequeno para terse producido razoablemente polo azar.

Problema 3.6. Un estudo dun ecosistema dun bosque de folla caduca indica que o promedio neto de transformacións de nitróxeno en nitrato presenta un incremento de 2Kg por hectárea e ano. Os enxeñeiros de montes cren que unha desfoliación da maleza do bosque conduciría a un descenso dese valor. Arráncase a maleza nun área de 15 hectáreas dun bosque experimental. Límpase a área para impedi-lo crecemento. Despois dun ano determinouse o cambio de nitróxeno a nitrato, por hectárea, analizando a auga da chuvia en 15 puntos dentro do bosque. Obtivéronse os seguintes resultados: $\bar{X} = -3$, $s_{n-1} = 7,5$. ¿Proba isto que arranca-la maleza do bosque provoca un descenso no cambio medio neto de nitróxeno a nitrato por hectárea e ano?

Solución. A variable aleatoria a considerar é X = “cambio neto de nitróxeno a nitrato por hectárea e ano”.

O contraste a considerar é

$$H_0: \mu \geq 2, \quad H_1: \mu < 2.$$

Temos como datos $\mu_0 = 2$, $n = 15$, $\bar{X} = -3$, $s_{n-1} = 7,5$.

Empregámo-lo estatístico de contraste $\frac{\bar{X}-\mu_0}{s_{n-1}/\sqrt{n}}$, que ten distribución t -Student.

Substituíndo,

$$\frac{-3 - 2}{7,5/\sqrt{15}} = -2,582.$$

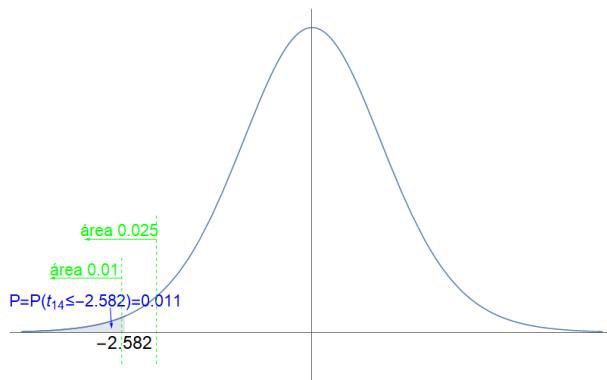


Figura 3.3: valor P

Resulta entón que o P -valor é $P = P(t_{14} \leq -2,582)$. Buscando nas táboas (hai que emprega-la simetría da t -Student e mira-lo valor á dereita de 2.582) obsérvase que $0,01 < P < 0,025$. Empregando software informático tense, de feito, $P = 0,011$.

Conclusión: como o valor P obtido é pequeno, *rexeitamos* H_0 e concluímos que hai evidencia significativa, polo menos do 97.5 %, de que a retirada de maleza do bosque deu como resultado un incremento inferior a 2Kg por hectárea e ano da concentración media de nitróxeno en forma de nitratos.

(Nótese que, se no enunciado do problema nos tivesen pedido un nivel de confianza do 99 %, teríamos que ter aceptado a hipótese nula, mentres que se o nivel de confianza fose do 97.5 % teríamos que tela rexeitado. Como o valor P se aproxima bastante ó 1 %, e cometer un erro de tipo I non parece que vaia ocasionar problemas graves, decidimos que o valor obtido é suficiente para rexeita-la hipótese nula.)

Cómpre enfatizar que como resultado da resolución deste problema, acabamos de probar que o *incremento de transformación de nitróxeno en nitrato é menor ca 2Kg*. Non estamos probando que diminúa a transformación de nitróxeno en nitrato (a pesar de que iso é o que pasa na mostra). Se quixeramos probar isto último, teríamos que face-lo contraste de hipóteses

$$H_0: \mu \geq 0, \quad H_1: \mu < 0.$$

Neste caso, o valor no estatístico resulta

$$\frac{-3 - 0}{7,5/\sqrt{15}} = -1,549,$$

e o valor P é $P = P(t_{14} < -1,549) = 0,0718$, que é un valor relativamente grande. Por tanto, para este problema teríamos que aceptar H_0 , e concluiríamos que non habería

evidencia significativa de que a as transformacións medias de nitróxeno en nitrato por hectárea e ano diminúisen. \square

Cando facemos contrastes bilaterais o procedemento máis común para un estatístico simétrico é considerar un valor P de dúas colas como dúas veces o valor P dunha cola. Non obstante, non existe consenso para calcula-lo valor nestes casos, especialmente se o estatístico non é simétrico.

3.2. Contraste de hipóteses para a varianza

Neste caso trátase de facer un contraste de hipótese sobre a varianza dunha poboación normal despois de ter elixido unha mostra aleatoria simple X_1, \dots, X_n .

3.2.1. Contrastes bilaterais

O contraste de hipóteses é neste caso

$$H_0: \sigma^2 = \sigma_0^2, \quad H_1: \sigma^2 \neq \sigma_0^2.$$

Suposto que H_0 fose certa, tomámo-lo estatístico de contraste

$$\frac{(n-1)s_{n-1}^2}{\sigma_0^2} \sim \chi_{n-1}^2,$$

que, como vimos na sección dedicada ó cálculo de intervalos de confianza para a varianza (Definición 2.18), ten unha distribución χ^2 con $n-1$ graos de liberdade (Definición 2.15).

Para un nivel de significación α temos:

- Rexión crítica: $[0, \chi_{1-\alpha/2}^2) \cup (\chi_{\alpha/2}^2, +\infty)$.
- Rexión de aceptación: $[\chi_{1-\alpha/2}^2, \chi_{\alpha/2}^2]$.

Se a media é coñecida empregáse o estatístico ns_{μ}^2/σ_0^2 e procédese de xeito análogo.

3.2.2. Contrastes unilaterais

Supoñemos agora que facemos un contraste unilateral dereito. Tamén suporemos que a media é descoñecida. Se non fose así tomaríamo-lo estatístico ns_{μ}^2/σ_0^2 e procederíamos similarmente. O contraste é entón

$$H_0: \sigma^2 \leq \sigma_0^2, \quad H_1: \sigma^2 > \sigma_0^2.$$

Tomámo-lo mesmo estatístico de contraste ca no caso anterior, co que para un nivel de significación α temos:

- Rexión crítica: $(\chi_{n-1, \alpha}^2, +\infty)$.
- Rexión de aceptación é $[0, \chi_{n-1, \alpha}^2]$.

Nótese que como esta distribución non é simétrica, para o contraste unilateral esquerdo haberá que tomar:

- Rexión crítica: $[0, \chi_{n-1, 1-\alpha}^2)$.
- Rexión de aceptación: $[\chi_{n-1, 1-\alpha}^2, +\infty)$.

Ó igual que sucedía co contraste de hipóteses unilateral para a media, unha alternativa para aceptar ou rexeita-la hipótese nula é calcula-lo valor P (Subsección 3.1.3) e comprobar se este número é pequeno ou non.

3.3. Contraste de hipóteses para unha proporción

Temos unha poboación na que unha determinada propiedade se dá con probabilidade p , tomamos unha mostra aleatoria simple X_1, \dots, X_n , e denotamos por \hat{p} á proporción desa propiedade que se dá na mostra.

3.3.1. Contrastes bilaterais

O contraste de hipóteses é neste caso

$$H_0: p = p_0, \quad H_1: p \neq p_0.$$

Suposto que H_0 fose certa, tomámo-lo estatístico de contraste

$$\frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \sim Z,$$

que, como vimos na sección dedicada ó cálculo de intervalos de confianza para a proporción (Definición 2.21), pode supoñerse que ten unha distribución normal $N(0, 1)$ se o tamaño da mostra n é suficientemente grande.

Para un nivel de significación α temos:

- Rexión crítica: $(-\infty, -Z_{\alpha/2}) \cup (Z_{\alpha/2}, +\infty)$, é dicir, cando $\left| \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \right| > Z_{\alpha/2}$.
- Rexión de aceptación: $[-Z_{\alpha/2}, Z_{\alpha/2}]$.

3.3.2. Contrastes unilaterais

Supoñemos agora que facemos un contraste unilateral dereito (para un contraste unilateral esquerdo procederíase de xeito análogo) da forma

$$H_0: p \leq p_0, \quad H_1: p > p_0.$$

Tomámo-lo mesmo estatístico de contraste ca no caso anterior, co que para un nivel de significación α temos:

- Rexión crítica: $(Z_{\alpha}, +\infty)$.

- Rexión de aceptación: $(-\infty, Z_\alpha]$.

Ó igual ca noutros casos, unha alternativa para aceptar ou rexeita-la hipótese nula é calcula-lo valor P (Subsección 3.1.3) e comprobar se este número é pequeno ou non.

Problema 3.7. Unha empresa farmacéutica quere comercializar un medicamento que cura certa doenza. Sábese que o 40 % dos doentes se curan sen toma-lo medicamento. A empresa debe probar que o seu medicamento é eficaz, e para iso adminístrao a 100 doentes, dos cales se curan 50. ¿É realmente eficaz o medicamento?

Solución. A variable aletoria a estudar é X , doentes que se curan despois de tomar certo medicamento.

A cuestión é se o medicamento cura máis ca non tomar nada. Para iso necesitarase evidencia concluínte de que na mostra se obtiveron resultados positivos. Por tanto, o contraste sobre proporcións é

$$H_0: p \leq 0,4, \quad H_1: p > 0,4.$$

Temos $p_0 = 0,4$, $n = 100$, e tomámo-lo estatístico de contraste $\frac{\hat{p}-p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$. Substituíndo os datos:

$$\frac{0,5 - 0,4}{\sqrt{\frac{0,4(1-0,4)}{100}}} = 2,04.$$

O valor P é por tanto $P = P(z \geq 2,04) = 0,0206$. De feito, $0,01 < P < 0,025$. En consecuencia, o resultado é significativo ó 2.5 %, pero non ó 1 %.

Conclusión: é dubidoso, pero como o nivel crítico é bastante pequeno, poderíamos rexeita-la hipótese nula e aceptar que existe evidencia significativa, polo menos do 97.5 %, de que a proporción de curacións entre as persoas que toman o medicamento é maior có 40 %. □

Observación 3.8. Unha cuestión que pode ser interesante é ternos preguntado, con anterioridade a ve-los resultados da mostraxe, polo número de casos que satisfán a propiedade buscada para que haxa que rexeita-la hipótese nula.

Poñamos, por exemplo, que témo-lo contraste de hipóteses unilateral esquerdo $H_0: p \geq p_0$, $H_1: p < p_0$. Supoñamos que o tamaño mostral é n , e que o nivel de significación é α . Temos que calcula-lo número máximo k de individuos para os que poderíamos rexeita-la hipótese nula H_0 .

Calculamos en primeiro lugar o valor Z_α . Así, para rexeitar H_0 , necesitamos que o valor no estatístico estea na rexión crítica $(-\infty, -Z_\alpha)$, é dicir,

$$\frac{\frac{k}{n} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} < -Z_\alpha.$$

Despexando k na anterior inecuación obtemos

$$k < n \left(p_0 - Z_\alpha \sqrt{\frac{p_0(1-p_0)}{n}} \right).$$

Problema 3.9. Para analiza-lo risco de sufrir un aborto espontáneo nos embarazos de mulleres hipertensas tratadas con inhibidores de enzima convertidora de angiotensina (IECA) durante o primeiro trimestre do embarazo, estudáronse 329 casos nos que se observaron 47 abortos espontáneos. Se a taxa de abortos espontáneos na poboación fose do 10 %,

1. ¿Poderíase afirmar que o tratamento con IECA no primeiro trimestre de embarazo incrementa a porcentaxe de abortos espontáneos?
2. ¿Cantos casos de abortos espontáneos terían que terse observado na mostra anterior para poder afirmar, cun nivel de significación do 0.05, que a taxa de abortos espontáneos en mulleres hipertensas sometidas a tratamento con IECA no primeiro trimestre de embarazo supera o 20 %?

Solución. Sexa X a variable aleatoria “abortos espontáneos en mulleres hipertensas sometidas a tratamento con IECA no primeiro trimestre de embarazo”.

Para a primeira parte do problema debemos face-lo contraste de proporcións

$$H_0: p \leq 0,1, \quad H_1: p > 0,1.$$

Os datos do problema dinnos que $p_0 = 0,1$, $n = 329$ e $\hat{p} = 47/329 = 0,143$. Substituíndo no estatístico de contraste obtemos

$$\frac{0,143 - 0,1}{\sqrt{\frac{0,1(1-0,1)}{329}}} = 2,59.$$

O valor P é $P = P(Z \geq 2,59) = 0,0048$, que é menor có 0.5 %.

Conclusión: *rexeitámo-la hipótese nula* e concluímos que hai evidencia significativa, polo menos do 99.5 %, de que o tratamento con IECA no primeiro trimestre de embarazo provoca que a porcentaxe de abortos espontáneos sexa maior có 10 %.

Para a segunda parte do problema témo-lo novo contraste de hipóteses

$$H_0: p \leq 0,2, \quad H_1: p > 0,2.$$

O nivel de significación é $\alpha = 0,05$. Así, $Z_\alpha = 1,6449$. Por tanto necesitamos atopar k na inecuación

$$\frac{\frac{k}{329} - 0,2}{\sqrt{\frac{0,2(1-0,2)}{329}}} > 1,6449.$$

Despexando obtemos $k > 77,73$.

Conclusión: necesitaríamos ter rexistrado polo menos 78 casos de abortos espontáneos nunha mostra de 329 mulleres para ter evidencia significativa, polo menos do 95 %, de que a taxa de abortos espontáneos en mulleres hipertensas sometidas a tratamento on IECA no primeiro trimestre de embarazo supera o 20 %. \square

3.4. Resumo de contrastes de hipóteses para unha poboación

A continuación preséntase unha táboa resumo cos resultados deste capítulo.

Contrastes para a media

Distribución na mostraxe	$\frac{\bar{X} - \mu_0}{s_{n-1}/\sqrt{n}} \sim t_{n-1}$
$H_0: \mu = \mu_0, H_1: \mu \neq \mu_0.$	
Rexión crítica	$(-\infty, -t_{n-1, \alpha/2}) \cup (t_{n-1, \alpha/2}, +\infty)$
Rexión de aceptación	$[-t_{n-1, \alpha/2}, t_{n-1, \alpha/2}]$
$H_0: \mu \leq \mu_0, H_1: \mu > \mu_0.$	
Rexión crítica	$(t_{n-1, \alpha}, +\infty)$
Rexión de aceptación	$(-\infty, t_{n-1, \alpha}]$

Contrastes para a varianza

Distribución na mostraxe	$\frac{(n-1)s_{n-1}^2}{\sigma_0^2} \sim \chi_{n-1}^2$
$H_0: \sigma^2 = \sigma_0^2, H_1: \sigma^2 \neq \sigma_0^2.$	
Rexión crítica	$[0, \chi_{n-1, 1-\alpha/2}^2) \cup (\chi_{n-1, \alpha/2}^2, +\infty)$
Rexión de aceptación	$[\chi_{n-1, 1-\alpha/2}^2, \chi_{n-1, \alpha/2}^2]$
$H_0: \sigma^2 \leq \sigma_0^2, H_1: \sigma^2 > \sigma_0^2.$	
Rexión crítica	$(\chi_{n-1, \alpha}^2, +\infty)$
Rexión de aceptación	$[0, \chi_{n-1, \alpha}^2]$
$H_0: \sigma^2 \geq \sigma_0^2, H_1: \sigma^2 < \sigma_0^2.$	
Rexión crítica	$[0, \chi_{n-1, 1-\alpha}^2)$
Rexión de aceptación	$(\chi_{n-1, 1-\alpha}^2, +\infty)$

Contrastes para a proporción

Distribución na mostraxe	$\frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}} \sim N(0, 1)$
$H_0: p = p_0, H_1: p \neq p_0.$	
Rexión crítica	$(-\infty, -Z_{\alpha/2}) \cup (Z_{\alpha/2}, +\infty)$
Rexión de aceptación	$[-Z_{\alpha/2}, Z_{\alpha/2}]$
$H_0: p \leq p_0, H_1: p > p_0.$	
Rexión crítica	$(Z_{\alpha}, +\infty)$
Rexión de aceptación	$(-\infty, Z_{\alpha}]$

Capítulo 4

Comparación de dúas poboacións

Neste capítulo centraremos no problema de comparar dúas poboacións. A situación xeral é que temos dúas poboacións de interese, e en ambas se trata de estudar a mesma característica. A cuestión é que as dúas poboacións se atopan, por así dicilo, en circunstancias distintas, e interesa comparalas para saber como afectan esas circunstancias particulares á medida da característica que se estuda. Colleremos unha mostra aleatoria simple en cada unha das poboacións, e a partir delas, trataremos de tomar unha decisión sobre a característica que estamos estudando.

En principio preséntanse dúas posibilidades para as mostras: que sexan independentes, ou que estean emparelladas.

Se as mostras son *independentes*, entón temos dúas poboacións para as cales estudamos respectivamente dúas variables aleatorias X e Y *independentes* cunhas distribucións de probabilidade que pertencen á mesma familia.

Extraemos unha mostra aleatoria simple X_1, \dots, X_{n_1} da primeira poboación, e Y_1, \dots, Y_{n_2} da segunda. Supoñemos que as dúas mostras son *independentes*, é dicir, que os obxectos ou individuos da mostra da primeira poboación non teñen relación algunha cos da segunda. Nótese que os tamaños mostrais n_1 e n_2 non teñen por que ser iguais.

Cando as mostras están emparelladas, o procedemento é distinto e tratarase máis adiante neste capítulo.

4.1. Comparación das medias de dúas poboacións con mostras independentes

Supoñamos que X e Y seguen as dúas unha distribución normal, $N(\mu_1, \sigma_1)$ e $N(\mu_2, \sigma_2)$, respectivamente. Un xeito de comparalas medias poboacionais é restalas e compara-la súa diferenza. (Aínda que as distribucións non sexan normais, se a mostra é suficientemente grande recordemos que podemos asumir os resultados que seguen en virtude do teorema central do límite.)

En primeiro lugar centrámonos na estimación puntual.

Definición 4.1. Se as medias mostrais son \bar{X} e \bar{Y} respectivamente, entón está claro que un estimador para a diferenza das medias é $\widehat{\mu_1 - \mu_2} = \bar{X} - \bar{Y}$.

Para o resto de consideracións desta sección temos varios casos.

4.1.1. Coñecidas as varianzas poboacionais

Definición 4.2. Se σ_1 e σ_2 son coñecidas, cousa que habitualmente non sucede, entón o estatístico

$$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim Z,$$

ten unha distribución normal estándar.

Coñecido tal estatístico podemos tanto calcular intervalos de confianza (seguindo o mesmo procedemento estudado no capítulo dedicado a intervalos de confianza (Capítulo 2)), como facer contrastes de hipótese (seguindo o procedemento do capítulo dedicado a contrastes de hipóteses (Capítulo 3)).

Intervalos de confianza

Observación 4.3. Por exemplo, un intervalo de confianza de nivel de significación α para a diferenza de medias vén determinado pola expresión

$$\left| \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \right| \leq Z_{\alpha/2}.$$

Por tanto, tal intervalo é da forma

$$\left[(\bar{X} - \bar{Y}) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, (\bar{X} - \bar{Y}) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right].$$

ou ben,

$$(\bar{X} - \bar{Y}) \pm Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}.$$

Contraste de hipóteses

Para un contraste bilateral

$$H_0: \mu_1 - \mu_2 = (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 \neq (\mu_1 - \mu_2)_0,$$

con nivel de significación α , temos:

- Rexión crítica: $(-\infty, -Z_{\alpha/2}) \cup (Z_{\alpha/2}, +\infty)$.
- Rexión de aceptación: $[-Z_{\alpha/2}, Z_{\alpha/2}]$.

Para un contraste unilateral

$$H_0: \mu_1 - \mu_2 \leq (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 > (\mu_1 - \mu_2)_0,$$

con nivel de significación α , temos:

- Rexión crítica: $(Z_{\alpha}, +\infty)$.
- Rexión de aceptación: $(-\infty, Z_{\alpha}]$.

Para un contraste unilateral esquerdo procederíase de xeito análogo.

4.1.2. Descoñecidas as varianzas poboacionais pero supostas iguais

Supoñamos agora que $\sigma^2 = \sigma_1^2 = \sigma_2^2$ é a varianza (coincidente) das dúas poboacións. Non obstante, σ é descoñecida.

En primeiro lugar temos que estima-la varianza. Para iso temos dous estimadores da mesma cantidade, $s_1^2 := s_{X, n_1-1}^2$ e $s_2^2 := s_{Y, n_2-1}^2$, obtidos a partir das dúas mostras. Para uni-la información obtida por ambos, calculámo-la cuasi-varianza mostral conxunta.

Definición 4.4. A *cuasi-varianza mostral conxunta* defínese como

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2},$$

que é a media ponderada das cuasi-varianzas das mostras.

Agora podemos considera-lo estatístico para resolve-los problemas.

Definición 4.5. O estatístico

$$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2},$$

ten unha distribución *t*-Student con $n_1 + n_2 - 2$ graos de liberdade.

Observación 4.6. Para determinar se podemos considera-la varianzas de dúas poboacións iguais ou non, unha posibilidade é empregar un test de hipóteses sobre a varianza, tal e como se describe na sección dedicada á comparación das varianzas de dúas poboacións con mostras independentes (Sección 4.2).

Outra posibilidade empregada habitualmente consiste en aceptar que $\sigma_1 = \sigma_2$ cando se ten

$$\frac{1}{2} \leq \frac{s_1^2}{s_2^2} \leq 2,$$

é dicir, cando ningunha das cuasi-varianzas é máis do dobre da outra.

Intervalos de confianza

Observación 4.7. Agora un intervalo de confianza de nivel de significación α para a diferenza de medias vén determinado pola expresión

$$\left| \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right| \leq t_{n_1+n_2-2, \alpha/2}.$$

Por tanto, tal intervalo é da forma (omitímo-los graos de liberdade)

$$\left[(\bar{X} - \bar{Y}) - t_{\alpha/2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, (\bar{X} - \bar{Y}) + t_{\alpha/2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right],$$

ou ben,

$$(\bar{X} - \bar{Y}) \pm t_{\alpha/2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}.$$

Contraste de hipóteses

Para un contraste bilateral

$$H_0: \mu_1 - \mu_2 = (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 \neq (\mu_1 - \mu_2)_0,$$

con nivel de significación α , temos:

- Rexión crítica: $(-\infty, -t_{\alpha/2}) \cup (t_{\alpha/2}, +\infty)$.
- Rexión de aceptación: $[-t_{\alpha/2}, t_{\alpha/2}]$.

Para un contraste unilateral

$$H_0: \mu_1 - \mu_2 \leq (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 > (\mu_1 - \mu_2)_0,$$

con nivel de significación α , temos:

- Rexión crítica: $(t_\alpha, +\infty)$.
- Rexión de aceptación: $(-\infty, t_\alpha]$.

En lugar de fixar un nivel de significación poderíamos ter calculado o valor P .

Para un contraste unilateral esquerdo procederíase de xeito análogo.

4.1.3. Descoñecidas as varianzas poboacionais

Se σ_1^2 e σ_2^2 son descoñecidas e non poden ser supostas iguais, entón non hai solución exacta para o problema de determina-la distribución na mostraxe de $\bar{X} - \bar{Y}$, o cal obriga a adoptar solucións aproximadas.

Cando os tamaños das mostraxes son grandes, cabe argumentar que s_1^2 e s_2^2 son boas aproximacións de σ_1^2 e σ_2^2 . Así,

$$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}},$$

terá unha distribución aproximadamente normal.

Non obstante, neste curso empregaremos-lo feito de que o estatístico anterior está mellor aproximado por unha t -Student cuns certos graos de liberdade.

Definición 4.8. A aproximación debida a Welch-Smith-Satterthwaite define a cantidade γ de graos de liberdade, como

$$\gamma \sim \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}}.$$

Como γ ten que ser enteiro, tomarase como valor a parte enteira do valor obtido no cálculo.

Definición 4.9. Así, o estatístico que empregaremos

$$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}},$$

terá, aproximadamente, unha distribución t de Student con γ graos de liberdade, onde γ está definido mediante a fórmula de Welch.

Intervalos de confianza

Observación 4.10. No caso máis xeral, un intervalo de confianza de nivel de significación α para a diferenza de medias vén determinado pola expresión

$$\left| \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \right| \leq t_{\gamma, \alpha/2}.$$

Por tanto, tal intervalo é da forma

$$\left[(\bar{X} - \bar{Y}) - t_{\gamma, \alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}, (\bar{X} - \bar{Y}) + t_{\gamma, \alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \right],$$

ou ben,

$$(\bar{X} - \bar{Y}) \pm t_{\gamma, \alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}.$$

Contraste de hipóteses

Para un contraste bilateral

$$H_0: \mu_1 - \mu_2 = (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 \neq (\mu_1 - \mu_2)_0,$$

con nivel de significación α , temos:

- Rexión crítica: $(-\infty, -t_{\gamma, \alpha/2}) \cup (t_{\gamma, \alpha/2}, +\infty)$.
- Rexión de aceptación: $[-t_{\gamma, \alpha/2}, t_{\gamma, \alpha/2}]$.

Para un contraste unilateral

$$H_0: \mu_1 - \mu_2 \leq (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 > (\mu_1 - \mu_2)_0,$$

con nivel de significación α , temos:

- Rexión crítica: $(t_{\gamma, \alpha}, +\infty)$.
- Rexión de aceptación: $(-\infty, t_{\gamma, \alpha}]$.

En troques de fixar un nivel de significación poderíamos ter calculado o valor P como

$$P = P(t_{\gamma} \geq \text{valor no estatístico}),$$

e decidir, se tal valor é moi pequeno, que rexeitámo-la hipótese nula; de non ser así, aceptámola.

Para un contraste unilateral esquerdo procederíase de xeito análogo.

Problema 4.11. Un isótopo radioactivo (Sr-90) acumúlase nos ósos por medio do leite de vaca consumido. Quérese coñecer se o nivel de isótopo nos nenos é distinto ca nos adultos. Para iso tómase:

- unha mostra aleatoria de 121 nenos; obtense unha concentración media de 2.6 picocurios/g, cunha cuasi-desviación típica de 1.2.
- outra mostra de 61 adultos; para estes tense como media 0.4 e cuasi-desviación típica 0.11.

Solución. Denotemos por X á variable aleatoria que mide o nivel de isótopo nos nenos, e por Y á dos adultos. Asumímo-las notacións que vimos empregando nesta sección. Trátase pois dun problema de contraste de hipóteses bilateral, é dicir,

$$H_0: \mu_1 = \mu_2, \quad H_1: \mu_1 \neq \mu_2.$$

Equivalentemente, miramos se a diferenza $\mu_1 - \mu_2$ é nula.

O problema dános como datos: $n_1 = 121$, $\bar{X} = 2,6$, $s_1 = 1,2$, $n_2 = 61$, $\bar{Y} = 0,4$, $s_2 = 0,11$.

Xa que non hai coñecemento das varianzas, e $s_1/s_2 = 1,2/0,11 = 10,9 > 2$, empregámo-la fórmula de Welch, co que substituíndo:

$$\gamma \sim \frac{\left(\frac{1,2^2}{121} + \frac{0,11^2}{61}\right)^2}{\frac{\left(\frac{1,2^2}{121}\right)^2}{120} + \frac{\left(\frac{0,11^2}{61}\right)^2}{60}} = 123,965,$$

Así que tomamos $\gamma = 123$. Realmente a distribución t_{123} é moi parecida á normal estándar.

Substituíndo no estatístico obtemos:

$$\frac{(2,6 - 0,4) - 0}{\sqrt{\frac{1,2^2}{121} + \frac{0,11^2}{61}}} = 20,001.$$

Este valor está fóra de intervalos da forma $[-t_{123, \alpha/2}, t_{123, \alpha/2}]$ para valores moito menores a $\alpha = 0,1\%$. (Por exemplo, $t_{123, 0,0005} = 3,371$.) O cálculo do valor P (Subsección 3.1.3) con software informático daría (nótese que é un contraste bilateral cun estatístico simétrico) $P = 2P(t_{123} > 20,001) = 0,2 \cdot 10^{-39}$, que é un valor moi pequeno.

Conclusión: *rexeitámo-la hipótese nula* e concluímos que hai evidencia significativa, cunha confianza de máis do 99.9%, de que o nivel de isótopo en nenos é distinto ó dos adultos. \square

Problema 4.12. En vista dos resultados obtidos no problema anterior, preguntámonos agora se hai evidencia significativa de que a media obtida para os nenos sexa maior cá dos adultos. ¿É maior ca 2 picocurios/g?

Solución. As variables aleatorias a considerar son as mesmas cá no problema anterior.

Como necesitamos evidencia concluínte de que a dos nenos é *significativamente* maior, temos que estuda-lo contraste:

$$H_0: \mu_1 \leq \mu_2, \quad H_1: \mu_1 > \mu_2.$$

Neste caso calculámo-lo valor P substituíndo no estatístico:

$$\begin{aligned} P &= P\left(t_{123} \geq \frac{(2,6 - 0,4) - 0}{\sqrt{\frac{1,2^2}{121} + \frac{0,4^2}{61}}}\right) \\ &= P(t_{123} \geq 20,001) = 0,9 \cdot 10^{-40}. \end{aligned}$$

Conclusión: rexeitámo-la hipótese nula e concluímos que hai evidencia significativa, cun nivel de confianza moi alto, de que o nivel de isótopo en nenos é superior ó dos adultos.

A última pregunta correspóndese a un contraste de hipóteses

$$H_0: \mu_1 - \mu_2 \leq 2, \quad H_1: \mu_1 - \mu_2 > 2.$$

Procédese igual, pero agora o valor no estatístico é

$$\frac{(2,6 - 0,4) - 2}{\sqrt{\frac{1,2^2}{121} + \frac{0,11^2}{61}}} = 1,818,$$

co cal, o valor P é $P = P(t_{123} \geq 1,818) = 0,0357$.

Agora $0,025 < P < 0,05$, así que parece que podemos rexeita-la hipótese nula con nivel de confianza do 95 %, pero non do 97.5 %.

Conclusión: rexeitámo-la hipótese nula e por tanto concluímos que hai evidencia significativa, polo menos ó 95 %, de que a diferenza de concentración de isótopo Sr-90 en nenos é superior a 2 picocurios/g con respecto á dos adultos. \square

4.2. Comparación das varianzas de dúas poboacións con mostras independentes

Na sección anterior vimos que é máis sinxelo, e que non hai que facer aproximacións, para estima-la diferenza de medias se supoñemos que as varianzas poboacionais son iguais (Subsección 4.1.2). A utilidade desta sección é precisamente dar un test de hipóteses para contrastar se as varianzas poboacionais son iguais.

De novo suporemos que X e Y seguen distribucións normais, $N(\mu_1, \sigma_1)$ e $N(\mu_2, \sigma_2)$, respectivamente. Extraemos dúas mostras independentes de cada poboación, X_1, \dots, X_{n_1} e Y_1, \dots, Y_{n_2} .

Como anteriormente, denotamos por s_1^2 e s_2^2 as cuasi-varianzas das mostras anteriores. Para comparar σ_1^2 e σ_2^2 , non é conveniente considera-lo estatístico $s_1^2 - s_2^2$. Por varias razóns é preferible emprega-lo cociente.

Definición 4.13. Para compara-las varianzas de dúas poboacións normais a partir de mostras independentes emprégase o estatístico

$$\frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} \sim F_{n_1-1, n_2-1},$$

que segue unha distribución F de **Fisher-Snedecor con $(n_1 - 1, n_2 - 1)$ graos de liberdade**.

Definición 4.14. A distribución de probabilidade da $F_{m,n}$ ten como función de densidade unha función da forma

$$f(x) = c_{m,n} x^{\frac{m}{2}-1} (n + mx)^{-\frac{m+n}{2}}, \quad x > 0,$$

sendo $c_{m,n}$ unha determinada constante.

A distribución F de Snedecor depende de dous parámetros, tamén chamados graos de liberdade, que haberá que considerar á hora de buscar valores nas táboas.

Proposición 4.15. *Algunhas propiedades da F de Snedecor:*

- $E(F_{m,n}) = \frac{n}{n-2}$ se $n > 2$.
- Só está definida para valores positivos e non é simétrica.
- Se $F_{m,n,\alpha}$ é o valor para o que $P(F_{m,n} > F_{m,n,\alpha}) = \alpha$, entón

$$F_{n,m,\alpha} = \frac{1}{F_{m,n,1-\alpha}}.$$

Observación 4.16. Nótese que nas táboas da F de Snedecor empregadas neste curso, o primeiro índice é o das columnas.

Contraste de hipóteses

Estamos interesados no contraste de hipóteses

$$H_0: \sigma_1^2 = \sigma_2^2, \quad H_1: \sigma_1^2 \neq \sigma_2^2.$$

Como por hipótese, $\sigma_1^2 = \sigma_2^2$, o estatístico anterior simplifícase e temos que empregar s_1^2/s_2^2 que, como vimos, ten distribución F_{n_1-1, n_2-1} . En consecuencia, para un nivel de significación α temos:

- Rexión crítica: $\left(0, \frac{1}{F_{n_2-1, n_1-1, \alpha/2}}\right) \cup \left(F_{n_1-1, n_2-1, \alpha/2}, +\infty\right)$.
- Rexión de aceptación: $\left[\frac{1}{F_{n_2-1, n_1-1, \alpha/2}}, F_{n_1-1, n_2-1, \alpha/2}\right]$.

Problema 4.17. Realizouse un estudo sobre as necesidades enerxéticas para o crecemento e mantemento dun niño de avións en Perthshire, Escocia. Obtivéronse os seguintes resultados para as observacións de dúas mostras independentes da variable normal “número de kilocalorías por gramo e hora que se requiren por paxaro”.

Adultos incubando	Adultos precriando
$n_1 = 57$	$n_2 = 12$
$\bar{X} = 0,0167$	$\bar{Y} = 0,0144$
$s_1 = 0,0042$	$s_2 = 0,0024$

¿Indican estes datos que o número de kilocalorías requerido por adultos que están incubando é maior có requerido polos adultos que están precriando? Razoalo empregando o valor P. (Facer un contraste de hipóteses para determina-la igualdade das varianzas empregando para iso $\alpha = 0,1$.)

Solución. Sexa X a variable aleatoria que mide o número de kilocalorías por gramo e hora que se requiren por paxaro en adultos incubando, e Y a que mide o mesmo valor en adultos precriando.

Trátase dun problema de contraste de hipóteses da media de dúas poboacións, que coa notación que vimos empregando, se escribe como

$$H_0: \mu_1 \leq \mu_2, \quad H_1: \mu_1 > \mu_2.$$

En primeiro lugar teremos que decidir se podemos supoñer que as varianzas poboacionais son ou non iguais. Isto require un contraste de hipóteses previo:

$$H_0: \sigma_1^2 = \sigma_2^2, \quad H_1: \sigma_1^2 \neq \sigma_2^2.$$

Substituíndo no estatístico s_1^2/s_2^2 obtemos $0,0042^2/0,0024^2 = 3,0625$. Por outro lado, para $\alpha = 0,1$ temos $F_{56,11,0,05} = 2,4960$ e $F_{56,11,0,95} = 0,5091$. Claramente $3,0625 \notin [0,5091, 2,4960]$, co que debemos rexeita-la última hipótese, e por tanto non podemos supoñer que σ_1 e σ_2 sexan iguais.

Temos que empregar entón a aproximación de Welch. Primeiro calculámo-los graos de liberdade:

$$\frac{\left(\frac{0,0042^2}{57} + \frac{0,0024^2}{12}\right)^2}{\frac{(0,0042^2/57)^2}{57-1} + \frac{(0,0024^2/12)^2}{12-1}} = 27,51,$$

así que tomamos $\gamma = 27$.

Substituímos agora no estatístico:

$$\frac{(0,0167 - 0,0144) - 0}{\sqrt{\frac{0,0042^2}{57} + \frac{0,0024^2}{12}}} = 2,589.$$

Entón o valor P é $P = P(t_{27} \geq 2,589) = 0,0077$, é dicir, $0,005 < P < 0,01$.

Conclusión: rexeitámo-la hipótese nula e concluímos que existe evidencia significativa, polo menos do 99%, de que o número de kilocalorías requerido por adultos incubando é maior có requerido polos adultos que están precriando. \square

Problema 4.18. Nun estudo sobre hábitos de alimentación de morcegos, márcanse 25 femias e 11 machos e rastréanse por radio. A variable de interese é “distancia que percorren voando en busca de alimento”. O experimento proporciona a seguinte información:

Femias	Machos
$n_1 = 25$	$n_2 = 11$
$\bar{X} = 205$	$\bar{Y} = 135$
$s_1 = 100$	$s_2 = 95$

Calcular un intervalo de confianza para a diferenza de medias cun nivel de confianza do 90%.

Solución. Sexa X a variable aleatoria que mide a distancia que percorre un morcego femia voando en busca de alimento, e Y a que mide a mesma distancia para morcegos macho.

Para calcula-lo intervalo de confianza pedido, primeiro temos que decidir que estatístico empregamos. Por tanto, hai que determinar se podemos supoñer que as varianzas poboacionais poden ser consideradas iguais ou non.

Así, investigámo-lo contraste de hipóteses

$$H_0: \sigma_1^2 = \sigma_2^2, \quad H_1: \sigma_1^2 \neq \sigma_2^2.$$

Substituíndo no estatístico s_1^2/s_2^2 obtemos $100^2/95^2 = 1,108$. Por outro lado, para $\alpha = 0,1$, $F_{24,10,0,05} = 2,737$ e $F_{24,10,0,95} = 0,4435$. Dado que o valor no estatístico cae entre os dous, aceptamos que as dúas varianzas poboacionais son iguais.

Podemos, por tanto, toma-la varianza conxunta

$$s_p^2 = \frac{(15 - 1) \cdot 100^2 + (11 - 1) \cdot 95^2}{25 + 11 - 2} = 9713,24.$$

Dado que $t_{34,0,05} = 1,691$ o intervalo de confianza buscado é

$$(205 - 135) \pm 1,691 \sqrt{9713,24} \sqrt{\frac{1}{25} + \frac{1}{11}} = 70 \pm 60,30,$$

que despois de face-los cálculos resulta $[9,70, 130,30]$.

Conclusión: cun nivel de confianza do 90 %, a diferenza entre a distancia media percorrida por un morcego femia e un morcego macho en busca de comida sitúase entre 9.70 e 130.30 metros.

Fixémonos que os resultados obtidos supoñendo que as varianzas poboacionais son iguais non é moi distinto do que se obtería se empregáramo-lo método xeral. Para o método xeral o número de graos de liberdade estímase por

$$\frac{\left(\frac{100^2}{25} + \frac{95^2}{11}\right)^2}{\frac{(100^2/25)^2}{25-1} + \frac{(95^2/11)^2}{11-1}} = 20,13,$$

de xeito que tomamos $\gamma = 20$.

A continuación calculamos $t_{20,0,05} = 1,72472$. Por tanto, o intervalo calcúlase como

$$(205 - 135) \pm 1,72472 \sqrt{\frac{100^2}{25} + \frac{95^2}{11}},$$

e facendo os cálculos resulta $[9,75, 130,25]$, que é moi similar ó obtido anteriormente. \square

4.3. Comparación de proporcións de dúas poboacións con mostras independentes

Supoñamos que hai dúas poboacións, nas que unha determinada propiedade se dá con probabilidades p_1 e p_2 , respectivamente. Tomamos dúas mostras aleatorias simples X_1, \dots, X_{n_1} e Y_1, \dots, Y_{n_2} de cada poboación.

De novo, o estimador puntual é o esperado:

Definición 4.19. Para estima-la diferenza de proporcións emprégase a diferenza das proporcións das mostras: $\widehat{p_1 - p_2} = \hat{p}_1 - \hat{p}_2$.

Definición 4.20. Tómase o estatístico

$$\frac{(\widehat{p_1 - p_2}) - (p_1 - p_2)}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}} \sim Z,$$

que ten distribución normal $N(0, 1)$.

4.3.1. Intervalos de confianza

Observación 4.21. Para un nivel de significación α , un intervalo de confianza para a diferenza de proporcions vén determinado pola expresión

$$\left| \frac{(\widehat{p_1 - p_2}) - (p_1 - p_2)}{\sqrt{\frac{\widehat{p_1}(1-\widehat{p_1})}{n_1} + \frac{\widehat{p_2}(1-\widehat{p_2})}{n_2}}} \right| \leq Z_{\alpha/2}.$$

Por tanto, un intervalo de confianza será da forma:

$$(\widehat{p_1 - p_2}) \pm Z_{\alpha/2} \sqrt{\frac{\widehat{p_1}(1-\widehat{p_1})}{n_1} + \frac{\widehat{p_2}(1-\widehat{p_2})}{n_2}}.$$

4.3.2. Contraste de hipóteses

Chamamos *valor nulo* ó valor fronte ó que contrastámo-la hipótese (o que en analogía con outros casos chamariamos $(p_1 - p_2)_0$). Dependendo de se este valor é cero ou non, podemos facer unha pequena simplificación. En particular, se o valor nulo é cero, non tomarémo-lo mesmo estatístico que para o cálculo dun intervalo de confianza.

Valor nulo distinto de cero

Para un contraste bilateral

$$H_0: p_1 - p_2 = (p_1 - p_2)_0, \quad H_1: p_1 - p_2 \neq (p_1 - p_2)_0,$$

con $(p_1 - p_2)_0 \neq 0$, e nivel de significación α , temos:

- Rexión crítica: $(-\infty, -Z_{\alpha/2}) \cup (Z_{\alpha/2}, +\infty)$.
- Rexión de aceptación: $[-Z_{\alpha/2}, Z_{\alpha/2}]$.

Para un contraste unilateral

$$H_0: p_1 - p_2 \leq (p_1 - p_2)_0, \quad H_1: p_1 - p_2 > (p_1 - p_2)_0,$$

con $(\widehat{p_1} - \widehat{p_2})_0 \neq 0$ e nivel de significación α , temos:

- Rexión crítica: $(Z_{\alpha}, +\infty)$.
- Rexión de aceptación: $(-\infty, Z_{\alpha}]$.

En troques de fixar un nivel de significación poderíamos calcula-lo valor P como temos feito en exemplos anteriores.

Para un contraste unilateral esquerdo procederíase de xeito análogo.

Valor nulo cero

Nun contraste bilateral escribimos

$$H_0: p_1 - p_2 = 0, \quad H_1: p_1 - p_2 \neq 0,$$

Como por hipótese as proporcións reais son as mesmas, \hat{p}_1 e \hat{p}_2 estiman a mesma cantidade. Así, facemos como que as dúas mostras proveñen dunha mesma poboación con proporción descoñecida p , e tomámo-la media ponderada das proporcións estimadas en cada mostra.

Definición 4.22. A media ponderada das proporcións para un contraste de hipóteses con valor nulo cero defínese como

$$\hat{p} = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}.$$

Definición 4.23. O estatístico de contraste neste caso simplifícase un pouco e tomamos

$$\frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \sim Z,$$

que tamén segue unha distribución normal estándar.

O resto das cuestións son análogas ó caso anterior.

Problema 4.24. Entre marzo e agosto de 1998 fíxose en Baltimore un ensaio clínico aleatorizado e dobre cego para comproba-la eficacia do paracetamol como analxésico para trata-la migraña. Un grupo de voluntarios da zona dividiuse aleatoriamente en dous grupos. A un deles subministróuselle paracetamol e ó outro un placebo. Entre os 147 pacientes que recibiron o paracetamol, 85 notaron diminución de dor ás dúas horas de tomalo, fronte a 56 de 142 no caso do grupo de control (o que tomou o placebo). A partir dos resultados deste estudo clínico, ¿hay evidencia suficiente, con nivel de confianza 99 %, para afirmar que o paracetamol é un analxésico útil á hora de trata-los síntomas da migraña? ¿É a diferenza de efectividade maior ó 10 %?

Solución. Sexa X a variable aleatoria que mide a eficacia do paracetamol como analxésico para trata-la migraña, Y a do placebo.

Trátase dun problema de contraste de hipóteses:

$$H_0: p_1 - p_2 \leq 0, \quad H_1: p_1 - p_2 > 0,$$

que por tanto ten valor nulo cero.

Témo-los datos: $n_1 = 147$, $\hat{p}_1 = 85/147 = 57,8\%$, $n_2 = 142$, $\hat{p}_2 = 56/142 = 39,4\%$. Así, a media ponderada das proporcións é:

$$\hat{p} = \frac{147 \cdot 0,578 + 142 \cdot 0,394}{147 + 142} = 0,488.$$

Substituíndo no estatístico de contraste:

$$\frac{0,578 - 0,394}{\sqrt{0,49481(1 - 0,49481)\left(\frac{1}{147} + \frac{1}{142}\right)}} = 3,126.$$

Por outra banda, $P = P(Z \geq 3,126) = 0,000886 < 0,01$ é un valor máis pequeno ca α .

Conclusión: *rexéitase a hipótese nula*, e concluímos que existe evidencia significativa, polo menos do 99 %, de que o paracetamol é útil no tratamento da migraña.

Con respecto á segunda pregunta, agora témo-lo contraste de hipóteses

$$H_0: p_1 - p_2 \leq 0,1, \quad H_1: p_1 - p_2 > 0,1.$$

Como o valor nulo non é cero, substituímos no estatístico xeral para obter

$$\frac{(0,578 - 0,394) - 0,1}{\sqrt{\frac{0,578(1-0,578)}{147} + \frac{0,394(1-0,394)}{142}}} = 1,4509$$

Agora $P = P(Z \geq 1,4509) = 0,0734 > 0,01$, é dicir, $P > 0,01$. Non podemos rexeita-la hipótese nula para o nivel de significación dado.

Conclusión: *aceptámo-la hipótese nula* e concluímos que non hai evidencia significativa de que a proporción de enfermos de migraña que melloran ás 2 horas de tomar paracetamol supere no 10 % ós que tomaron o placebo.

En consecuencia, hai evidencias de que, ante un ataque de migraña, é mellor tomar paracetamol que non tomar nada. Non obstante, a diferenza de proporcións entre enfermos que toman paracetamol e que non toman nada, e melloran ás 2 horas, non supera o 10 %.

Observación 4.25. Determinación do tamaño da mostra

Podemos facer un argumento similar a cando estimámo-lo tamaño da mostra para unha proporción (Observación 2.25), e así, se facemos $n \leq n_1, n_2$, para un erro máximo ϵ obtemos (vendo que o máximo de $x(1-x)$ está en $x = 1/2$)

$$Z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \leq Z_{\alpha/2} \sqrt{\frac{1}{4n} + \frac{1}{4n}} \leq \epsilon.$$

Despexando, $n \geq \frac{Z_{\alpha/2}^2}{2\epsilon^2}$. En consecuencia, teremos que tomar

$$n_1, n_2 \geq \frac{Z_{\alpha/2}^2}{2\epsilon^2}.$$

4.4. Comparación da media con mostras emparelladas

Nesta sección estudamos un caso que se presenta con bastante frecuencia. É aquel no que as dúas mostras están emparelladas, é dicir, que para cada individuo dunha lle corresponde un da outra, asociado de xeito natural ou a propósito. Isto dáse por exemplo cando se estuda o comportamento de dous xemelgos, nais e fillas, ou o comportamento dunha persoa antes e despois de tomar un medicamento. Nestes casos faise un único sorteo e a segunda mostra dedúcese da primeira.

Os métodos das seccións anteriores *non* son aplicables xa que, para calcula-los estatísticos correspondentes, facíase uso da independencia entre mostras.

Temos, por tanto, dúas poboacións nas que medimos unha mesma característica, e denotamos por X e Y as variables aleatorias de cada unha. Tomamos X_1, \dots, X_n unha mostra aleatoria simple, que está emparellada con Y_1, \dots, Y_n , non independente da anterior, e do mesmo tamaño. Considerámo-la variable aleatoria $D = X - Y$, que supoñemos que está normalmente distribuída. Así, temos unha mostra das diferencias D_1, \dots, D_n , onde $D_i = X_i - Y_i$.

En consecuencia, acabamos de reduci-lo problema de dúas mostras a facer inferencia sobre a súa diferenza. No caso da diferenza de medias temos que estudar entón $\mu_D = \mu_1 - \mu_2$. Como vimos, o estatístico necesario para estudar-la media é

$$\frac{\bar{D} - \mu_D}{s_D/\sqrt{n}} \sim t_{n-1},$$

que segue unha distribución t -Student con $n - 1$ graos de liberdade.

O procedemento para obter intervalos de confianza e facer contrastes de hipóteses é, por tanto, o mesmo có discutido para o cálculo de intervalos de confianza para a media (Subsección 2.2.3), e contraste de hipóteses para a media (Sección 3.1), respectivamente.

Observación 4.26. Cabe resaltar que, se ben $\mu_D = \mu_1 - \mu_2$ e $\bar{D} = \bar{X} - \bar{Y}$, pola contra $s_D^2 \neq s_1^2 \pm s_2^2$. En consecuencia, a varianza da diferenza non pode ser deducida das varianzas das dúas variables.

Problema 4.27. Realízase un estudo para investiga-lo efecto do exercicio no nivel de colesterol no sangue (mg/dl). Tomáronse mostras de sangue nos participantes. Despois someuse ós individuos a un programa de exercicios, e volvéronse tomar mostras de sangue. Obtivéronse os seguintes datos:

Persoa	Nivel previo	Nivel posterior
1	182	198
2	232	210
3	191	194
4	200	220
5	148	138
6	249	220
7	276	219
8	213	161
9	241	210
10	480	313
11	262	226

Construír un intervalo de confianza para a diferenza de medias con nivel de confianza do 90%.

Solución. En primeiro lugar organizámo-los cálculos para a diferenza, o cadrado das diferencias, e sumámo-los resultados. Denotamos por X á variable aleatoria que mide o nivel de colesterol en sangue antes do exercicio, e Y o nivel despois do exercicio. Tomámo-la diferenza $D = X - Y$.

Persoa	X	Y	D	D^2
1	182	198	-16	256
2	232	210	22	484
3	191	194	-3	9
4	200	220	-20	400
5	148	138	10	100
6	249	220	29	841
7	276	219	57	3249
8	213	161	52	2704
9	241	210	31	961
10	480	313	167	27889
11	262	226	36	1296
Σ	2674	2309	365	38189

En vista da táboa temos

$$\bar{D} = \frac{1}{n} \sum_i D_i = \frac{365}{11} = 33,182,$$

$$s_D^2 = \frac{1}{n-1} \sum_i D_i^2 - \frac{n}{n-1} \bar{D}^2 = \frac{38189}{10} - \frac{11}{10} 33,182^2 = 2607,76.$$

Ademais, $t_{10,0,05} = 1,812$. Así, o intervalo de confianza vén dado por

$$33,18 \pm 1,81 \frac{\sqrt{2607,76}}{\sqrt{11}},$$

o que, facendo os cálculos resulta [5,28, 61,09].

Conclusión: cun nivel de confianza do 90 %, a diferenza media de nivel de colesterol entre a xente que fai exercicio e a que non sitúase entre 5.28 e 61.09. \square

4.5. Resumo de contrastes de hipóteses para dúas poboacións

A continuación preséntase unha táboa resumo cos resultados deste capítulo ata agora.

Comparación da media de dúas poboacións	
Varianzas poboacionais descoñecidas pero iguais	
Estatístico	$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2}$
Cuasi-varianza conxunta	$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$
<i>Intervalos de confianza</i>	
Inecuación	$\left \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right \leq t_{n_1+n_2-2, \alpha/2}$
Fórmula	$(\bar{X} - \bar{Y}) \pm t_{n_1+n_2-2, \alpha/2} s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$
$H_0: \mu_1 - \mu_2 = (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 \neq (\mu_1 - \mu_2)_0.$	
Rexión crítica	$(-\infty, -t_{n_1+n_2-2, \alpha/2}) \cup (t_{n_1+n_2-2, \alpha/2}, +\infty)$
$H_0: \mu_1 - \mu_2 \leq (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 > (\mu_1 - \mu_2)_0.$	
Rexión crítica	$(t_{n_1+n_2-2, \alpha}, +\infty)$
Varianzas poboacionais descoñecidas	
Estatístico	$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim t_\gamma$
Graos de liberdade	$\gamma \sim \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}}$
<i>Intervalos de confianza</i>	
Inecuación	$\left \frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \right \leq t_{\gamma, \alpha/2}$
Fórmula	$(\bar{X} - \bar{Y}) \pm t_{\gamma, \alpha/2} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$
$H_0: \mu_1 - \mu_2 = (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 \neq (\mu_1 - \mu_2)_0.$	
Rexión crítica	$(-\infty, -t_{\gamma, \alpha/2}) \cup (t_{\gamma, \alpha/2}, +\infty)$
$H_0: \mu_1 - \mu_2 \leq (\mu_1 - \mu_2)_0, \quad H_1: \mu_1 - \mu_2 > (\mu_1 - \mu_2)_0.$	
Rexión crítica	$(t_{\gamma, \alpha}, +\infty)$
Comparación da varianza de dúas poboacións	
Estatístico	$\frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} \sim F_{n_1-1, n_2-1}$
$H_0: \sigma_1^2 = \sigma_2^2, \quad H_1: \sigma_1^2 \neq \sigma_2^2.$	
Rexión crítica	$\left(0, \frac{1}{F_{n_2-1, n_1-1, \alpha/2}}\right) \cup (F_{n_1-1, n_2-1, \alpha/2}, +\infty)$

Comparación de proporcións de dúas poboacións	
Estatístico	$\frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}} \sim Z$
<i>Intervalos de confianza</i>	
Inecuación	$\left \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}} \right \leq Z_{\alpha/2}$
Fórmula	$(\hat{p}_1 - \hat{p}_2) \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$
$H_0: p_1 - p_2 = (p_1 - p_2)_0 \neq 0, \quad H_1: p_1 - p_2 \neq (p_1 - p_2)_0.$	
Rexión crítica	$(-\infty, -Z_{\alpha/2}) \cup (Z_{\alpha/2}, +\infty)$
$H_0: p_1 - p_2 \leq (p_1 - p_2)_0 \neq 0, \quad H_1: p_1 - p_2 > (p_1 - p_2)_0.$	
Rexión crítica	$(Z_{\alpha}, +\infty)$
Valor nulo cero	
Estatístico	$\frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}} \sim Z$
Proporción ponderada	$\hat{p} = \frac{n_1\hat{p}_1 + n_2\hat{p}_2}{n_1 + n_2}$
$H_0: p_1 = p_2, \quad H_1: p_1 \neq p_2.$	
Rexión crítica	$(-\infty, -Z_{\alpha/2}) \cup (Z_{\alpha/2}, +\infty)$
$H_0: p_1 \leq p_2, \quad H_1: p_1 > p_2.$	
Rexión crítica	$(Z_{\alpha}, +\infty)$

Capítulo 5

A proba chi-cadrado

A proba chi-cadrado, ou proba χ^2 , é un contraste de hipóteses introducido por Pearson para determinar se a discrepancia entre as frecuencias esperadas e as frecuencias observadas nunha táboa de continxencia é estatisticamente significativa.

Neste capítulo, as variables estatísticas son discretas: só toman un número finito de valores, divididos en categorías. Distinguiremos dous tipos de tests, que computacionalmente son practicamente iguais, pero que conceptualmente son un pouco distintos.

5.1. Contrastes de independencia para datos categóricos

Supoñamos que nunha poboación estamos interesados en observar dúas características X e Y que se corresponden con datos categóricos, é dicir, que son datos nominais que soamente poden tomar valores concretos, chamados categorías. Cada valor está nunha, e só nunha, categoría (é dicir, as categorías son disxuntas). Poñamos que X pode tomar f valores distintos A_1, \dots, A_f , e que Y pode tomar c valores B_1, \dots, B_c . O problema ó que nos enfrentamos agora é o de determinar se as dúas características X e Y son ou non independentes. De feito, o que queremos é ver se hai (ou non) evidencia significativa de que as dúas características non son independentes.

Tomamos unha mostra aleatoria simple bidimensional (é dicir, medindo as dúas características) na poboación, $(X_1, Y_1), \dots, (X_n, Y_n)$. Denotamos por n_{ij} ó número de observacións na mostra de tal xeito que o valor de X se atopa en A_i e o valor de Y en B_j . Os valores poden por tanto dispoñerse nunha *táboa de continxencia*, que consiste en organiza-los datos do seguinte xeito:

$X \setminus Y$	B_1	B_2	...	B_c	Σ
A_1	n_{11}	n_{12}	...	n_{1c}	$n_{1\cdot}$
A_2	n_{21}	n_{22}	...	n_{2c}	$n_{2\cdot}$
\vdots	\vdots	\vdots		\vdots	\vdots
A_f	n_{f1}	n_{f2}	...	n_{fc}	$n_{f\cdot}$
Σ	$n_{\cdot 1}$	$n_{\cdot 2}$...	$n_{\cdot c}$	n

Nesta táboa empregouse a notación:

$$n_{i\cdot} = n_{i1} + n_{i2} + \cdots + n_{ic},$$

$$n_{\cdot j} = n_{1j} + n_{2j} + \cdots + n_{fj},$$

que son os números totais de observacións que se atopan nos conxuntos A_i e B_j respectivamente. Obviamente, $n_{1\cdot} + \cdots + n_{f\cdot} = n_{\cdot 1} + \cdots + n_{\cdot c} = n$.

As probabilidades reais da poboación son denotadas como

$$p_{ij} = P((X \in A_i) \cap (Y \in B_j)),$$

e así, por se-las categorías disxuntas,

$$p_{i\cdot} = P(X \in A_i) = p_{i1} + p_{i2} + \cdots + p_{ic},$$

$$p_{\cdot j} = P(Y \in B_j) = p_{1j} + p_{2j} + \cdots + p_{fj}.$$

Isto podería organizarse tamén nunha táboa de probabilidades:

$X \setminus Y$	B_1	B_2	...	B_c	Σ
A_1	p_{11}	p_{12}	...	p_{1c}	$p_{1\cdot}$
A_2	p_{21}	p_{22}	...	p_{2c}	$p_{2\cdot}$
\vdots	\vdots	\vdots		\vdots	\vdots
A_f	p_{f1}	p_{f2}	...	p_{fc}	$p_{f\cdot}$
Σ	$p_{\cdot 1}$	$p_{\cdot 2}$...	$p_{\cdot c}$	1

En caso de que as dúas características fosen realmente independentes, teriamos

$$p_{ij} = P((X \in A_i) \cap (Y \in B_j)) = P(X \in A_i)P(Y \in B_j) = p_{i\cdot} p_{\cdot j}$$

para calquera para i, j . En consecuencia, o contraste de hipóteses que pretendemos estudar é:

$$H_0: p_{ij} = p_{i\cdot} p_{\cdot j}, \quad \forall i \in \{1, \dots, f\}, \forall j \in \{1, \dots, c\}.$$

Por tanto, o que resta por facer é estima-las probabilidades p_{ij} e atopar un estatístico convinte que nos permita decidir se cos valores da mostra podemos ou non descartar H_0 .

A partir da mostra, os estimadores obvios das probabilidades son

$$\widehat{p}_{ij} = \frac{n_{ij}}{n}, \quad \widehat{p}_{i\cdot} = \frac{n_{i\cdot}}{n}, \quad \widehat{p}_{\cdot j} = \frac{n_{\cdot j}}{n}.$$

Por outra banda, baixo a hipótese de independencia, o valor da celda (i, j) da táboa de continxencia debería ser

$$E_{ij} = np_{ij} = n p_{i\cdot} p_{\cdot j},$$

que por tanto se estima por

$$\widehat{E}_{ij} = n \widehat{p}_{i\cdot} \widehat{p}_{\cdot j} = \frac{n_{i\cdot} n_{\cdot j}}{n}.$$

Para determinar se os n_{ij} están suficientemente próximos a \widehat{E}_{ij} , empregámo-lo estatístico

$$\sum_{i,j} \frac{(n_{ij} - \widehat{E}_{ij})^2}{\widehat{E}_{ij}} \sim \chi_{(f-1)(c-1)}^2,$$

que segue aproximadamente unha distribución χ^2 de Pearson con $(f-1)(c-1)$ graos de liberdade cando a mostra é suficientemente grande.

Este contraste é unilateral dereito, así que para un nivel de significación α , temos

- Rexión crítica: $(\chi^2_{(f-1)(c-1), \alpha}, \infty)$.
- Rexión de aceptación: $[0, \chi^2_{(f-1)(c-1), \alpha}]$.

Obviamente, tamén se podería calcula-lo valor P e rexeita-la hipótese nula cando este valor sexa moi pequeno.

Problema 5.1. Realízase un estudo para investiga-la asociación entre a cor e a fragancia das azaleas silvestres. Obsérvanse 200 prantas floridas seleccionadas aleatoriamente, e clasifícase cada unha delas segundo a cor e a presenza de fragancia.

fragancia \ cor	branca	rosa	naraxa
si	12	60	58
non	50	10	10

¿Hai probas significativas de asociación entre a cor das flores e a súa fragancia?

Solución. Denotemos por X a fragancia dunha azalea, e por Y a súa cor. En primeiro lugar construímo-la táboa de continxencia:

fragancia \ cor	branca	rosa	naraxa	Σ
si	12	60	58	130
non	50	10	10	70
Σ	62	70	68	200

O problema consiste en facer un contraste de hipóteses de independencia para datos categóricos, é dicir,

$$H_0: p_{ij} = p_i \cdot p_j, \quad \forall i \in \{1, 2\}, \quad \forall j \in \{1, 2, 3\}.$$

Veremos máis adiante que a razón de que este sexa un contraste de independencia é que o investigador simplemente clasifica os datos do total da mostra en dúas categorías (neste caso, fragancia e cor das azaleas).

Para resolve-lo problema, calculámo-los valores esperados, no suposto de que houberse independencia das variables, mediante a fórmula $\widehat{E}_{ij} = \frac{n_i \cdot n_j}{n}$ (en verde), e tamén os valores $(n_{ij} - \widehat{E}_{ij})^2 / \widehat{E}_{ij}$ (en vermello), obtendo:

fragancia \ cor	branca	rosa	naraxa	Σ
si	12 40.3 19.87	60 45.5 4.62	58 44.2 4.31	130
non	50 21.7 36.91	10 24.5 8.58	10 23.8 8.00	70
Σ	62	70	68	200

Finalmente aprovéitanse todos estes cálculos para determina-lo valor no estatístico, (que consiste en suma-los valores vermellos), para obter 82.29.

O estatístico segue unha distribución χ^2 con $(2 - 1)(3 - 1) = 2$ graos de liberdade. Xa que non nos dan un nivel de significación, calculámo-lo valor P como $P = P(\chi^2_2 \geq 82,29) < 0,001$. (Utilizando software informático obtense $P = 1,35 \cdot 10^{-18}$.)

Conclusión: rexeitámo-la hipótese nula, e concluimos que si hai evidencia significativa, cun nivel de confianza moi alto (maior có 99.9%), de que existe relación entre a cor da flor e a súa fragancia. □

5.2. Contrastes de homoxeneidade para datos categóricos

Este contraste é bastante parecido ó da sección anterior, polo menos no que a cálculos se refire, anque o obxectivo é bastante distinto.

Supoñamos que temos f poboacións nas que se observa unha determinada característica que pode tomar un valor de entre c valores distintos A_1, \dots, A_c . O problema ó que nos enfrentamos é o de determinar se a distribución de probabilidade desa característica é a mesma en todas esas poboacións, ou se polo contrario, ditas poboacións son heteroxéneas con distintas distribucións de probabilidade.

Tomamos unha mostra aleatoria simple en cada unha das poboacións, con tamaños n_1, \dots, n_f , respectivamente. Denotamos por n_{ij} o número de observacións na mostra i que se atopa en A_j . Os datos poden dispoñerse nunha *táboa de continxencia*, organizada do seguinte xeito:

Mostra	A_1	A_2	...	A_c	tamaño
1	n_{11}	n_{12}	...	n_{1c}	n_1
2	n_{21}	n_{22}	...	n_{2c}	n_2
\vdots	\vdots	\vdots		\vdots	\vdots
f	n_{f1}	n_{f2}	...	n_{fc}	n_f
Σ	$n_{\cdot 1}$	$n_{\cdot 2}$...	$n_{\cdot c}$	n

De novo empregouse a notación:

$$n_{\cdot j} = n_{1j} + n_{2j} + \dots + n_{fj},$$

que son os números totais de observacións que se atopan nos conxuntos A_i . Ademais, $n = n_1 + \dots + n_f$ é o tamaño que se obtén ó xuntar tódalas mostras.

A hipótese de homoxeneidade significa que cada conxunto A_j ten unha probabilidade p_j independente da poboación i . Por tanto, se p_{ij} é a probabilidade de A_j na poboación i , a hipótese nula é

$$H_0: p_{1j} = p_{2j} = \dots = p_{fj} (= p_j), \quad \forall j \in \{1, \dots, c\}.$$

As probabilidades p_j poden estimarse mediante

$$\widehat{p}_j = \frac{n_{\cdot j}}{n}.$$

Baixo a hipótese de homoxeneidade, a frecuencia teórica de A_j na poboación i é

$$E_{ij} = n_i p_j,$$

que por tanto se estima mediante

$$\widehat{E}_{ij} = \frac{n_i n_{\cdot j}}{n}.$$

Para determinar se os n_{ij} están suficientemente próximos a \widehat{E}_{ij} empregámo-lo estatístico

$$\sum_{i,j} \frac{(n_{ij} - \widehat{E}_{ij})^2}{\widehat{E}_{ij}} \sim \chi_{(f-1)(c-1)}^2,$$

que segue aproximadamente unha distribución χ^2 de Pearson con $(f - 1)(c - 1)$ graos de liberdade.

Este contraste é unilateral dereito, así que para un nivel de significación α , temos

- Rexión crítica: $(\chi^2_{(f-1)(c-1), \alpha}, \infty)$.
- Rexión de aceptación: $[0, \chi^2_{(f-1)(c-1), \alpha}]$.

O contraste de independencia e o contraste de homoxeneidade son moi similares en canto a cálculos e interpretación. A diferenza fundamental está no xeito de selecciona-las mostras, xa que no contraste de homoxeneidade, o tamaño das mostras (é dicir, o total das filas) está fixado polo experimentador, mentres que no contraste de independencia é arbitrario.

Problema 5.2. Para probar unha nova vacina contra a hepatite, tómanse 549 voluntarios ós que se lles administra a vacina, e 534 ós que non. Ó cabo dun tempo obsérvanse os seguinte casos de enfermidade:

mostra	ten hepatite	tamaño
vacinado	11	549
non vacunado	70	534

¿É a vacina eficaz?

Solución. Para ver se a vacina é eficaz temos que dar evidencia significativa de que a proporción de enfermos de hepatite é menor na poboación dos vacunados. Por tanto, é un contraste de hipóteses sobre homoxeneidade no que pretendemos refuta-la hipótese nula de que a distribución de probabilidade do número de enfermos de hepatite é a mesma para as dúas poboacións.

En primeiro lugar completámo-la táboa de continxencia:

	hepatite si	hepatite non	tamaño
vacinado	11	538	549
non vacunado	70	464	534
Σ	81	1002	1083

Temos que facer un contraste de hipóteses de homoxeneidade para datos categóricos:

$$H_0: p_{11} = p_{21}, p_{12} = p_{22}.$$

A continuación calculámo-los valores esperados no suposto de que houbose homoxeneidade nas poboacións mediante a fórmula $\widehat{E}_{ij} = \frac{n_i n_{.j}}{n}$ (en verde), e tamén os valores $(n_{ij} - \widehat{E}_{ij})^2 / \widehat{E}_{ij}$ (en vermello), obtendo:

vacinado \ hepatite	si	non	tamaño
si	11 41.06 22.01	538 507.94 1.78	549
non	70 39.94 22.63	464 494.06 1.83	534
Σ	81	1002	1083

Finalmente aprovéitanse todas estas contas para calcula-lo valor no estatístico, (que consiste en suma-los valores vermellos), para obter 48.24.

O estatístico segue unha distribución χ^2 con $(2 - 1)(2 - 1) = 1$ grao de liberdade. Xa que non nos dan un nivel de significación, calculámo-lo valor P como $P = P(\chi_1^2 \geq 48,24) < 0,001$. (Empregando software informático obtense $P = 3,7 \cdot 10^{-12}$.)

Conclusión: rexeitámo-la hipótese nula, e concluimos que hai evidencia significativa, cun nivel de confianza superior ó 99.9%, de que a proporción de enfermos de hepatite é distinta dependendo de se estamos na poboación de individuos vacinados ou non vacinados.

En realidade, este contraste de hipóteses non serve para determinar se o medicamento é eficaz ou non. Non obstante, tendo en conta os datos da táboa, onde se observa que a proporción de enfermos de hepatite na poboación dos individuos vacinados é menor cá esperada, podemos concluir que a vacina é eficaz. \square

Os contrastes de independencia e homoxeneidade son bastante populares e empréganse a miúdo como estudos preliminares para ver se hai relación entre dúas ou máis variables. Nótese non obstante, que estes contrastes non nos din *cal* é a relación entre as variables, aínda que mirando os valores da táboa podemos sacar algunha conclusión. Para atopar unha relación que explique como se relaciona unha variable con outra necesitanse outras técnicas estatísticas como a regresión.

É moi típico que, por erro, descoñecemento, ou por tratar de influencia-la opinión da xente, se extraían dun test deste estilo conclusións distintas (aínda que aparentemente relacionadas) ás que en realidade se fan no estudo. Tamén é típico extrapolar causalidade (un suceso implica outro), cando só hai correlación (dous sucesos pasan ó mesmo tempo).

Problema 5.3. Realízase un estudo de mercado consistente en clasifica-la poboación de acordo co seu poder adquisitivo en nivel alto, medio e baixo. Tómase unha mostra de 50 persoas de cada clase social e mírase se posúen un reloxo de marca Rolex. Constátase que da clase alta teñen un 30 persoas, 14 de clase media, e 5 de clase baixa.

- Realiza-lo correspondente contraste de hipóteses para ver se existe relación entre a clase social e ser posuidor dun Rolex.
- ¿Podemos afirmar que hai evidencia estatística de que mercar un Rolex aumenta o poder adquisitivo?

Solución. Temos 3 poboacións, dependendo do “poder adquisitivo”, e a variable aleatoria Y =“ter un Rolex”.

En primeiro lugar construímo-la táboa de continxencia:

renta \ Rolex	si	non	tamaño
alta	30	20	50
media	14	36	50
baixa	5	45	50
Σ	49	101	150

Temos que face-lo contraste de hipóteses:

$$H_0: p_{11} = p_{21} = p_{31}, p_{12} = p_{22} = p_{32}.$$

Este é un contraste de hipóteses para homoxeneidade de datos categóricos (Sección 5.2), xa que o tamaño da mostra en cada poboación é fixado polo investigador. Para iso empregámo-lo estatístico

$$\sum_{i,j} \frac{(n_{ij} - \widehat{E}_{ij})^2}{\widehat{E}_{ij}},$$

que segue unha distribución χ^2 de Pearson con $(f-1)(c-1)$ graos de liberdade.

O número de graos de liberdade da distribución é $(3-1)(2-1) = 2$.

A continuación calculámo-las frecuencias esperadas, no suposto de que a hipótese nula sexa certa, mediante a fórmula $\widehat{E}_{ij} = \frac{n_{i.}n_{.j}}{n}$:

renta \ Rolex	si	non	tamaño
alta	16.33	33.67	50
media	16.33	33.67	50
baixa	16.33	33.67	50
Σ	49	101	150

Agora calculámo-los valores intermedios do estatístico $(n_{ij} - \widehat{E}_{ij})^2 / \widehat{E}_{ij}$:

renta \ Rolex	si	non	Σ
alta	11.435	5.548	
media	0.333	0.162	
baixa	7.864	3.815	
Σ			29.157

A suma dos valores intermedios, que coincide co valor no estatístico, é 29.157.

Calculámo-lo valor P como $P = P(\chi_2^2 > 29,157) = 0,5 \cdot 10^{-6}$, que é un valor pequeno.

Conclusión: Rexeitamos H_0 , e concluímos que hai evidencia significativa, cun nivel de confianza moi elevado, de que ter un Rolex depende do poder adquisitivo da persoa.

Non obstante, os datos que se dan neste exercicio non están encamiñados a responde-la segunda pregunta. Podemos deducir que hai relación entre “ter poder adquisitivo alto” e “ter un Rolex”, pero non podemos dicir nada con respecto á relación entre “mercar un Rolex” e “aumenta-lo poder adquisitivo”. Aínda que a lóxica indica a pensar que a resposta á segunda pregunta é negativa, os datos do problema non o confirman nin o refutan. \square

5.3. Bondade do axuste

Ata agora os contrastes de hipóteses foron empregados para decidi-la veracidade dunha hipótese sobre os parámetros dunha distribución. En ocasións, non obstante, é necesario emitir un xuízo sobre a distribución poboacional no seu conxunto. O problema da bondade do axuste consiste en decidir, á vista dos datos dunha mostra aleatoria simple dunha poboación, se pode admitirse que a distribución poboacional coincide cunha certa distribución dada (no noso caso unha $N(0, 1)$). Nótese que este é un problema *non paramétrico*.

Supoñamos que queremos averiguar se a distribución F dunha poboación se axusta a unha distribución normal $N(\mu, \sigma)$. Supoñemos que temos unha mostra aleatoria simple

X_1, \dots, X_n . O noso contraste é por tanto

$$H_0: F = N(\mu, \sigma), \quad H_1: F \neq N(\mu, \sigma).$$

En primeiro lugar teremos que estima-los valores dos parámetros. Empregaremos para iso os estimadores insesgados $\hat{\mu} = \bar{X}$ e $\hat{\sigma} = s_{n-1}$.

O segundo paso deste contraste consiste en agrupa-los datos en intervalos. Para iso realízase o seguinte procedemento:

- Busca-los valores máis pequeno e máis grande. Tomar valores “redondos” un pouco máis pequenos có máis pequeno, e un pouco máis grande có máis grande tendo en conta a precisión dos datos.
- Decidir cantos intervalos se van empregar. Dividiranse os datos en intervalos do mesmo tamaño (preferentemente os extremos deberían ser enteiros, ou números “redondos”). Hai varias regras para decidir este número. Unha posibilidade é tomar aproximadamente \sqrt{n} intervalos. É conveniente que cada intervalo resultante teña polo menos 5 valores. O número de tales intervalos denotámolo por r .
- Calcula-los límites dos intervalos tendo en conta os datos anteriores.
- Face-lo recuento de valores en cada intervalo.

Unha vez que témo-los datos divididos en intervalos, podemos calcula-la frecuencia observada o_i destes en cada intervalo. Este datos compáranse coa probabilidade de que a distribución $N(\mu, \sigma)$ estea entre cada un dos valores dos intervalos, multiplicada por n . Isto é o que se chama a frecuencia teórica e_i . Utilízase para iso a estimación de μ e σ .

Para decidir se as discrepancias entre as frecuencias mostrais e as teóricas son significativas, emprégase a proba χ^2 de Pearson. Tomamos por tanto o estatístico

$$\chi_{r-k-1}^2 = \sum_{i=1}^r \frac{(o_i - e_i)^2}{e_i},$$

que segue unha distribución χ^2 de Pearson con $r - k - 1$ graos de liberdade, onde k é o número de parámetros que tivemos que estimar para precisa-la distribución teórica (se son μ e σ , entón $k = 2$).

O estatístico anterior úsase para facer un *contraste unilateral dereito*.

Problema 5.4. Unha máquina produce pezas cunha determinada lonxitude, a cal se quere saber se segue unha distribución normal. Obtense a seguinte mostra:

10.39	10.66	10.12	10.32	10.25	10.91	10.52	10.83	10.72	10.28
10.35	10.46	10.54	10.72	10.23	10.18	10.62	10.49	10.32	10.61
10.64	10.23	10.29	10.78	10.81	10.39	10.34	10.62	10.75	10.34
10.41	10.81	10.64	10.53	10.31	10.46	10.47	10.43	10.57	10.74

Deséxase saber se esta mostra avala a hipótese de que a máquina produce pezas cunha lonxitude que efectivamente é normal.

Solución. Sexa X a variable aleatoria “lonxitude das pezas que produce a máquina”.

Vemos que temos $n = 40$ datos. En primeiro lugar estimamos puntualmente a media e a desviación típica. Para iso obtemos $\bar{X} = 10,502$ e $s_{n-1} = 0,205$.

Trátase por tanto do contraste de hipóteses

$$H_0: F = N(\bar{X}, s_{n-1}), \quad H_1: F \neq N(\bar{X}, s_{n-1}),$$

pero agora non contrastamos ou estimámo-lo valor dos parámetros, se non o feito de que a distribución sexa ou non normal.

Para face-lo contraste de χ^2 de bondade de axuste, primeiramente temos que dividi-lo percorrido dos valores en intervalos. Como hai 40 datos, dividimos en 7 intervalos, o que é aproximadamente $\sqrt{40}$. O mínimo é 10,12 e o máximo 10,91. Podemos tomar como rango $[10, 11]$ e dividilo en 7. Isto dá $(11 - 10)/7 = 0,1428$, así que redondeamos a 0,15 e repartímo-lo exceso $0,15 \cdot 7 - 1 = 0,05$ a cada lado. Así, os intervalos serían:

$$\begin{array}{lll} (9.975, 10.125] & (10.125, 10.275] & (10.275, 10.425] \\ (10.425, 10.575] & (10.575, 10.725] & (10.725, 10.875] \\ (10.875, 11.025] & & \end{array}$$

A continuación contámo-lo número de elementos en cada subintervalo:

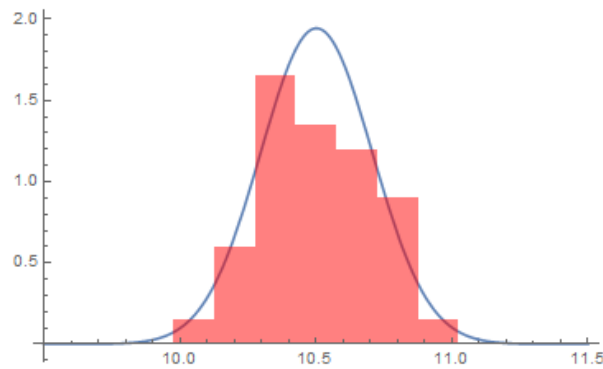


Figura 5.1: Gráfico de barras de frecuencias

Intervalo	o_i
(9.975, 10.125]	1
(10.125, 10.275]	4
(10.275, 10.425]	11
(10.425, 10.575]	9
(10.575, 10.725]	8
(10.725, 10.875]	6
(10.875, 11.025]	1

Temos $r = 7$. Ademais, co propósito de comparar coa distribución teórica, que é unha normal $N(10,502, 0,205)$, temos que toma-las colas ata $-\infty$ e $+\infty$. Completámo-la última columna cos valores teóricos e_i , que corresponden coa probabilidade de que a distribución estea no intervalo, multiplicada por n .

Intervalo	o_i	e_i
$(-\infty, 10,125]$	1	1.3223
$(10,125, 10,275]$	4	4.04803
$(10,275, 10,425]$	11	8.77801
$(10,425, 10,575]$	9	11.4123
$(10,575, 10,725]$	8	8.89851
$(10,725, 10,875]$	6	4.16001
$(10,875, \infty]$	1	1.38084
Σ	40	40.

Unha vez temos tódolos datos, só queda substituír no estatístico de contraste con $k = 2$, $r = 7$. Este estatístico segue unha distribución χ^2 con $r - k - 1 = 4$ graos de liberdade. Así

$$\sum_{i=1}^7 \frac{(o_i - e_i)^2}{e_i} = 2,16109,$$

e obtemos $P(\chi_4^2 \geq 2,16109) = 0,706158$, que é un valor moi alto.

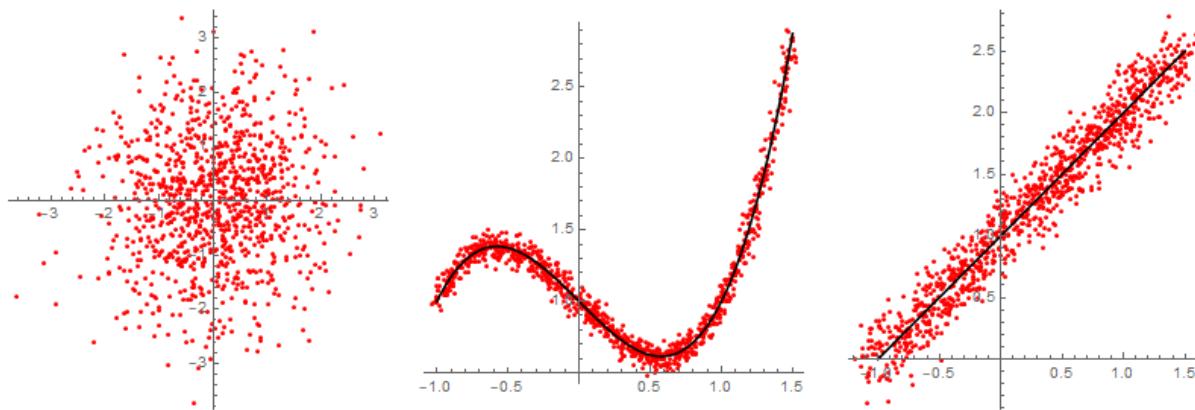
Conclusión: aceptámo-la hipótese de que a lonxitude das pezas da máquina segue unha distribución normal. \square

Capítulo 6

Regresión e correlación

O obxectivo deste capítulo é tratar de establecer a dependencia dunhas variables aleatorias con outras. En principio asumiremos que un determinado efecto se pode explicar mediante unhas causas e un erro. Asumiremos que temos dúas variables aleatorias X e Y . O obxectivo é atopar unha función f tal que $Y = f(X) + \epsilon$. Así, a Y chámase *resposta*, a f a *explicación*, e ϵ é o *erro*.

O seguinte gráfico amosa tres nubes de puntos distintas obtidas despois de tomar unha mostra aleatoria de dúas variables X e Y . No primeiro caso é evidente que non existe moita relación entre as dúas variables. No segundo caso parece que as variables están bastante relacionadas, e salvo un pequeno erro, dá a impresión de que Y se explica como dependente de X a través dunha ecuación polinómica. Finalmente, a terceira nube de puntos semella que se axusta a unha recta, aínda que o erro cometido é considerablemente máis grande ca no segundo exemplo.



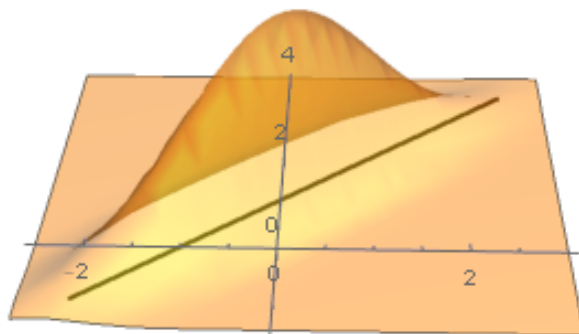
Nos dous últimos debuxos anteriores, é claro que existe unha dependencia (máis ou menos forte) entre X e Y . O obxectivo dun modelo de regresión é:

- Coñecer de que xeito a variable Y depende de X . Isto é o que se chama construír un *modelo de regresión*.
- Unha vez construído o modelo de regresión, empregar este para determina-lo valor de Y cando o valor de X é coñecido.

Neste capítulo consideraremos soamente o caso do *modelo de regresión linear simple*, que é aquel no que as variables X e Y son unidimensionais (como habitualmente), e que

Y se explica a partir de X mediante a ecuación dunha recta (coma no terceiro debuxo). Tamén se tratarán outros modelos que se reducen facilmente do de regresión linear.

6.1. Regresión linear



Sexan X e Y dúas variables aleatorias. O modelo de regresión linear consiste en atopar a recta $y = \alpha + \beta x$ que minimiza

$$E[(Y - (\alpha + \beta X))^2],$$

onde o que se trata é de atopar α e β . Non é difícil ver que estes dous valores se poden calcular simplemente derivando a anterior expresión con respecto de α e de β e igualando a cero. (Analogamente a como se fai para calculalo mínimo dunha función, pero con dúas variables.) Esta recta chámase a recta de regresión mínimo-cuadrática, porque na práctica se obtén despois de minimizala distancia cuadrática media dos puntos dunha mostra a dita recta.

Despois de face-los cálculos resulta que a ecuación da recta buscada é

$$Y - \mu_Y = \frac{\sigma_{XY}}{\sigma_X^2}(X - \mu_X) + \epsilon,$$

onde

- $\mu_X = E[X]$ é a media de X ,
- $\mu_Y = E[Y]$ é a media de Y ,
- $\sigma_X^2 = E[(X - \mu_X)^2]$ é a varianza de X ,
- $\sigma_Y^2 = E[(Y - \mu_Y)^2]$ é a varianza de Y ,
- $\sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)] = E[XY] - E[X]E[Y]$ é a **covarianza** entre X e Y ,
- ϵ é unha variable aleatoria que representa o erro cometido.

É consecuencia da construción do modelo que o erro ten media cero $\mu_\epsilon = E[\epsilon] = 0$, e que a súa varianza $\sigma_\epsilon^2 = V[\epsilon]$ é mínima.

Defínese o **coeficiente de correlación** de Pearson coma o cociente

$$\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y},$$

que satisfai $-1 \leq \rho \leq 1$, e dá unha idea do bo que é o axuste.

6.1.1. Estimación dos valores

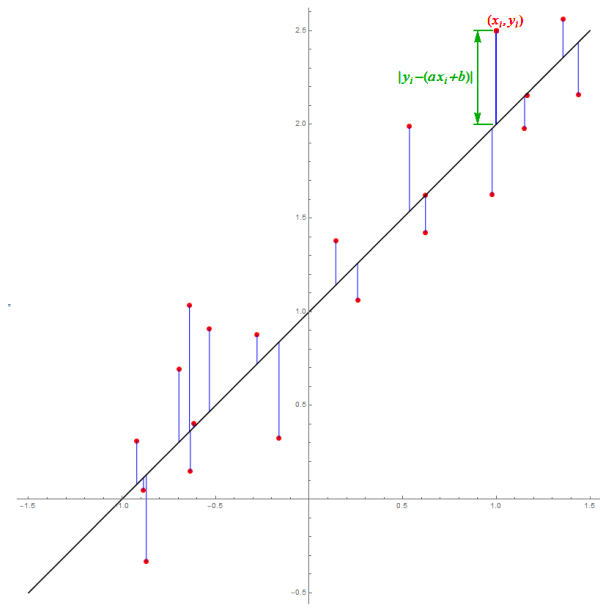


Figura 6.1: Distancia vertical entre a recta de regresión e un punto da mostra

Na práctica as variables aleatorias X e Y non son coñecidas e son estimadas por valores concretos $(x_1, y_1), \dots, (x_n, y_n)$ dunha mostra. Por iso, o modelo de regresión linear estímase como

$$Y - \bar{y} = \frac{s_{XY}}{s_X^2}(X - \bar{x}) + \epsilon,$$

onde agora

$$\begin{aligned}\bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i, \\ \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i, \\ s_X^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2, \\ s_Y^2 &= \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^n y_i^2 - \bar{y}^2, \\ s_{XY} &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}.\end{aligned}$$

Unha estimación equivalente para a recta de regresión é:

$$Y = a + bX + \epsilon,$$

onde

$$\begin{aligned}b &= \hat{\beta} = \frac{s_{XY}}{s_X^2}, \\ a &= \hat{\alpha} = \bar{y} - b\bar{x}.\end{aligned}$$

Nótese que esta recta de regresión sempre pasa por (\bar{x}, \bar{y}) .

6.1.2. Covarianza e correlación

A covarianza é a forma máis común de medi-la relación linear entre dúas variables. Para datos concretos recordemos que se estima por

$$\begin{aligned}s_{XY} &= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}.\end{aligned}$$

A covarianza non se ve afectada por cambios de posición, pero si de escala. De feito,

$$s_{aX+b, cY+d} = ac s_{XY}.$$

Para obter unha medida da relación linear entre dúas variables que non dependa da escala introduciuse o coeficiente de correlación, que para datos concretos se estima mediante

$$r = \frac{s_{XY}}{s_X s_Y}.$$

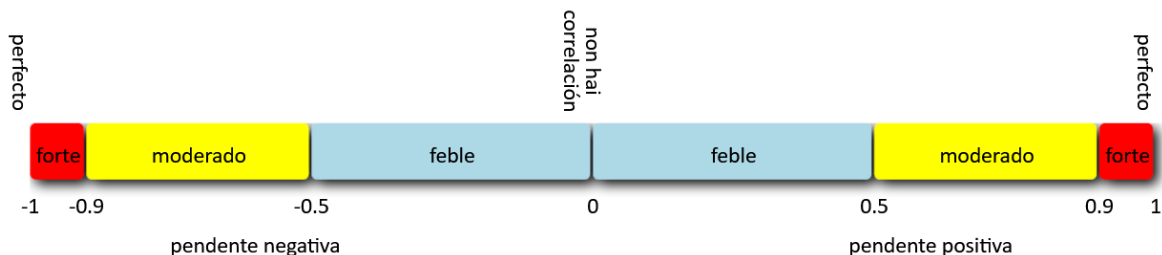
Unha versión equivalente, pero máis estable numericamente, é

$$r = \frac{n \sum_{i=1}^n x_i y_i - \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right)}{\sqrt{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \sqrt{n \sum_{i=1}^n y_i^2 - \left(\sum_{i=1}^n y_i \right)^2}}.$$

Proposición 6.1. *O coeficiente de correlación satisfai as seguintes propiedades:*

- *O coeficiente de correlación ten o mesmo signo cá pendente da recta de regresión.*
- *$-1 \leq r \leq 1$; valores próximos a 0 indican que o axuste é malo, valores próximos a 1 indican que o axuste é bo e que a relación é crecente, mentres que valores próximos a -1 indican que o axuste é bo e que a relación é decrecente.*

Un rango de valores para a bondade do axuste en función de r pode se-lo seguinte:



Algúns exemplos de coeficientes de correlación poden verse na táboa adxunta (Figura 6.2).

Problema 6.2. Os seguintes datos correspóndense co tempo transcorrido e a velocidade de caída dun obxecto:

tempo	velocidade
1	20.52
2	29.14
3	36.76
4	47.80
5	58.72

Calcula-la recta de regresión e dar unha aproximación da aceleración da gravidade. ¿Como de bo é o axuste?

Solución. Temos dúas variables que chamaremos t (tempo) e v (velocidade). En primeiro lugar dispoñémo-los cálculos:

	t	v	t^2	tv	v^2
1.	1.	20.52	1.	20.52	421.07
2.	2.	29.14	4.	58.28	849.14
3.	3.	36.76	9.	110.28	1351.30
4.	4.	47.80	16.	191.20	2284.84
5.	5.	58.72	25.	293.60	3448.04
Σ	15.	192.94	55.	673.88	8354.39

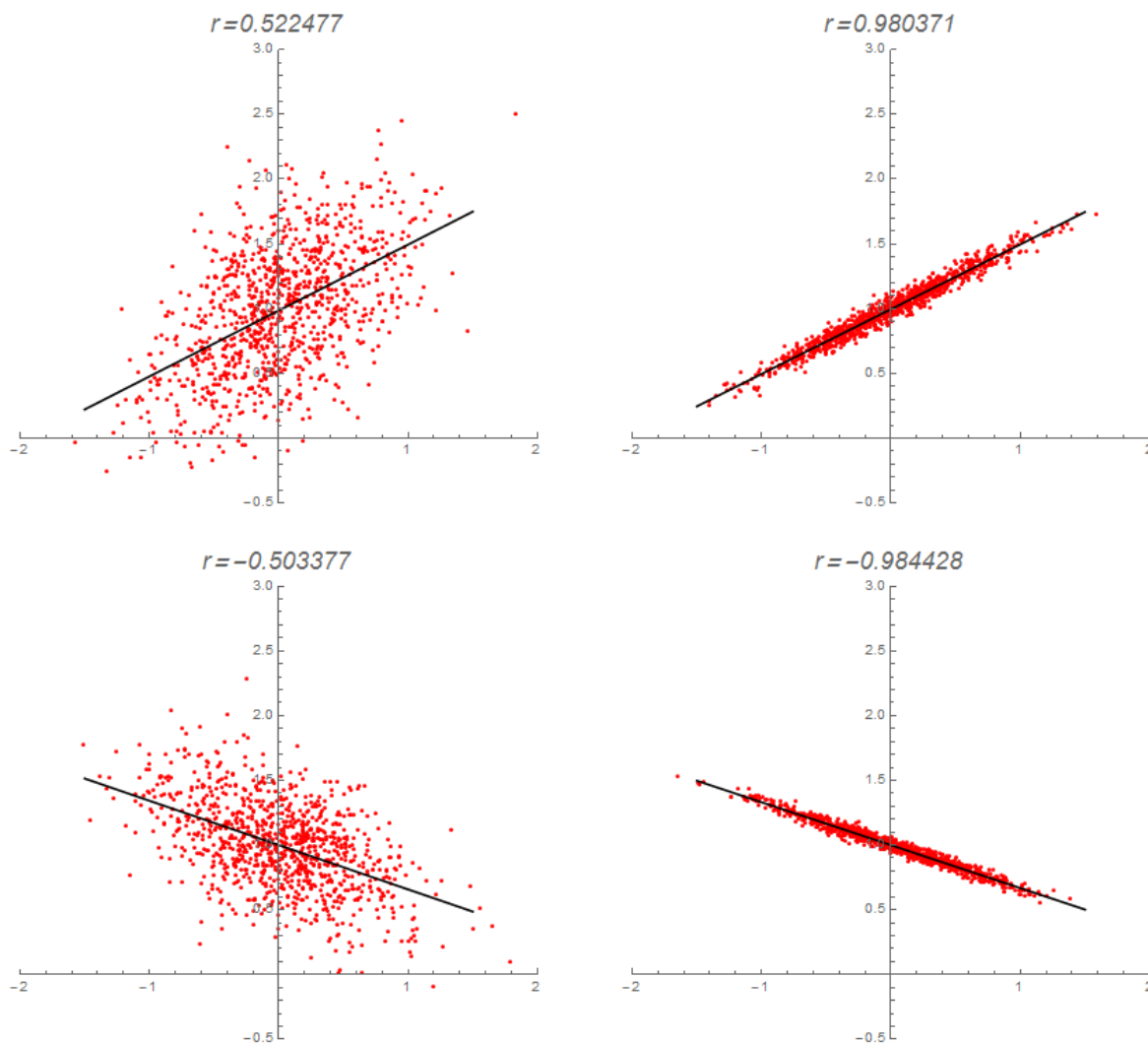


Figura 6.2: Coeficiente de correlación e calidade do axuste

Entón

$$\begin{aligned}\bar{t} &= \frac{15}{5} = 3,0, \\ \bar{v} &= \frac{192,94}{5} = 38,59, \\ s_t^2 &= \frac{55}{5} - 3^2 = 2,0, \\ s_v^2 &= \frac{8354,39}{5} - 38,59^2 = 181,84, \\ s_{tv} &= \frac{673,88}{5} - 3 \cdot 38,59 = 19,01,\end{aligned}$$

co que, substituíndo na fórmula, obtémo-la recta de regresión $v - 38,59 = \frac{19,01}{2}(t - 3)$, ou ben, como $b = 19,01/2,0 = 9,51$ e $a = 38,59 - 9,51 \cdot 3,0 = 10,07$, que

$$v = 10,07 + 9,51t,$$

de onde ademais se deduce que, en vista do resultado coñecido de física $v = v_0 + gt$, que

$g = 9,51$ é unha aproximación da aceleración da gravidade.

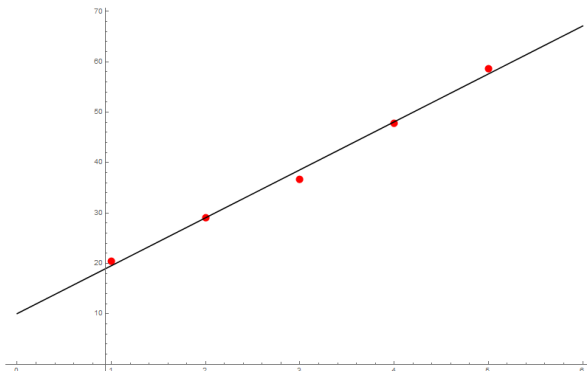


Figura 6.3: Os puntos e a súa recta de regresión

Finalmente calculámo-lo coeficiente de correlación:

$$r = \frac{19,01}{\sqrt{2}\sqrt{181,84}} = 0,997,$$

o cal quere dicir que o axuste é bo. □

6.1.3. Regresión exponencial

O procedemento para calcular unha regresión linear pode ser empregado tamén noutros contextos simplemente facendo un pequeno cambio de variable. Por exemplo, supoñamos que temos dúas variables aleatorias Z e T , e cremos que Z se explica a partir de T a través dunha fórmula exponencial:

$$Z = z_0 e^{-kT},$$

onde z_0 e k son os parámetros que queremos determinar. Entón, tomando logaritmos (neperianos)

$$\log Z = \log(z_0 e^{-kT}) = \log z_0 - kT.$$

Chamando $Y = \log Z$, $X = T$, $b = -k$, $a = \log z_0$, estamos exactamente na situación $Y = a + bX$ do principio. Por tanto, este tipo de axuste exponencial redúcese a un axuste linear, que xa sabemos resolver.

Problema 6.3. Inxectamos por vía intravenosa $125mg$ dun medicamento. Témo-las seguintes concentracións plasmáticas a medida que pasa o tempo:

tempo	concentración
1	5.0
2	3.0
3	2.0
4	1.5

Queremos estima-la curva exponencial da concentración de medicamento en sangue.

Solución. É sabido que a evolución da concentración teórica C dun medicamento en sangue ó longo do tempo t segue unha curva exponencial $C = c_0 e^{-kt}$. Despois de tomar logaritmos neperianos temos $\log C = \log c_0 - kt$, así que para calcula-la recta de regresión destes datos organizámo-los cálculos do seguinte xeito:

	$X = t$	C	$Y = \log C$	X^2	XY	Y^2
	1.	5.0	1.61	1.	1.61	1.59
	2.	3.0	1.10	4.0	2.20	1.21
	3.	2.0	0.69	9.0	2.08	0.48
	4.	1.5	0.41	16.0	1.62	0.16
Σ	10.	11.5	3.81	30.0	7.51	4.44

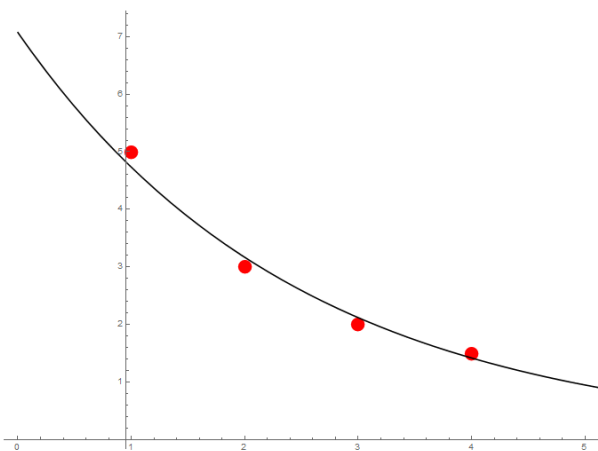


Figura 6.4: Os puntos e a súa regresión exponencial

Entón

$$\begin{aligned}\bar{X} &= \frac{10}{4} = 2,5, \\ \bar{Y} &= \frac{3,81}{4} = 0,95, \\ s_X^2 &= \frac{30,0}{4} - 2,5^2 = 1,25, \\ s_Y^2 &= \frac{4,44}{4} - 0,95^2 = 0,20, \\ s_{XY} &= \frac{7,51}{4} - 2,5 \cdot 0,95 = -0,50,\end{aligned}$$

co que, substituíndo na fórmula, obtémo-la recta de regresión:

$$Y - 0,95 = -\frac{0,50}{1,25}(X - 2,5).$$

Equivalentemente, obtense $b = -0,50/1,25 = -0,40$, $a = 0,95 + 0,40 \cdot 2,5 = 1,96$, de onde se deduce $Y = 1,96 - 0,40X$, ou $\log C = 1,96 - 0,40t$. Desfacendo o cambio de variable obtemos

$$C = 7,07e^{-0,40t}.$$

Pódese ver ademais que o coeficiente de correlación é

$$r = \frac{-0,50}{\sqrt{1,25}\sqrt{0,20}} = -0,992,$$

o que, ademais dun bo axuste, indica que a variable Y (ou a concentración C) decrece en función do tempo. \square

6.1.4. Regresión potencial

A regresión potencial é un caso bastante parecido ó da regresión exponencial. Neste caso hai dúas variables P e A que están relacionadas mediante a fórmula

$$P = \alpha A^\beta.$$

Para resolver isto, tomamos coma na sección anterior logaritmos e obtemos

$$\log P = \log(\alpha A^\beta) = \log \alpha + \beta \log A.$$

Así, chamando $Y = \log P$ e $X = \log A$ volvemos estar nun caso de axuste linear, que xa vimos como se resolve.

6.2. Análise da varianza

O obxectivo desta sección é estudar con máis profundidade se o modelo de regresión construído é correcto e útil. Para iso imos empregar un método coñecido como ANOVA (analysis of variance).

En primeiro lugar recordamos que $Y = \alpha + \beta X + \epsilon$, onde $\hat{Y} = \alpha + \beta X$ será a estimación dada pola recta de regresión, e $\epsilon = Y - \hat{Y}$ é o erro. Un cálculo non trivial amosa que as varianzas están relacionadas mediante

$$\sigma_Y^2 = \sigma_{\hat{Y}}^2 + \sigma_\epsilon^2.$$

Isto significa que a *variabilidade da variable dependente* Y , σ_Y^2 , se descompón como

- A *variabilidade explicada*, $\sigma_{\hat{Y}}^2$, que é aquela que se pode explicar en base ó modelo de regresión. De feito, como $\hat{Y} = \alpha + \beta X$, entón $\sigma_{\hat{Y}}^2 = \beta^2 \sigma_X^2$.
- A *variabilidade residual*, σ_ϵ^2 , que é a que non explica o modelo de regresión.

Chámase **coeficiente de determinación** á proporción entre a variabilidade explicada e a variabilidade da variable dependente. Por tanto,

$$\frac{\sigma_{\hat{Y}}^2}{\sigma_Y^2} = \frac{\beta^2 \sigma_X^2}{\sigma_Y^2} = \frac{\sigma_{XY}^2}{\sigma_X^2 \sigma_Y^2} = \rho^2,$$

que é o cadrado do coeficiente de correlación. Por tanto, $0 \leq \rho^2 \leq 1$.

Como xa sucedía co coeficiente de correlación, se $\rho^2 = 1$ (é dicir, se $\rho = \pm 1$) entón o axuste é perfecto. Valores de ρ^2 próximos a 1 significan que o axuste é bo, mentres que valores próximos a 0 indican un axuste malo.

Ademais, das fórmulas anteriores temos que

$$\begin{aligned} \sigma_{\hat{Y}}^2 &= \rho^2 \sigma_Y^2, \\ \sigma_\epsilon^2 &= \sigma_Y^2 - \sigma_{\hat{Y}}^2 = \sigma_Y^2 - \rho^2 \sigma_Y^2 = (1 - \rho^2) \sigma_Y^2. \end{aligned}$$

6.2.1. ANOVA

En realidade, os cálculos da sección anterior son teóricos, porque en xeral as distribucións X e Y non son coñecidas. Na práctica tómanse unha mostra e utilízanse as estimacións escritas con anterioridade.

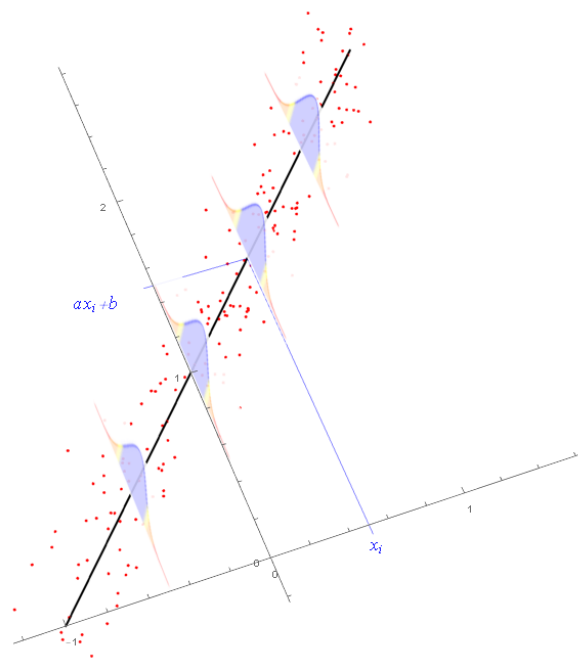


Figura 6.5: $Y \mid X = x_i$ está normalmente distribuída

Para continuar supoñamos que estamos traballando con n valores específicos x_1, \dots, x_n . Por tanto, os valores da variable explicativa están fixados polo experimentador e non son aleatorios. Só é aleatorio o erro, e en consecuencia a variable resposta. Unha mostra resultante deste tipo de experimento (chamado de deseño fixo), é do tipo $(x_1, Y_1), \dots, (x_n, Y_n)$. Asumimos que as variables aleatorias $Y \mid X = x_1, \dots, Y \mid X = x_n$ seguen distribucións normais independentes coa mesma varianza σ^2 . Se a regresión linear é válida, as medias destas variables están xustamente en $a + bx_i$, é dicir, $(Y \mid X = x_i) \sim N(a + bx_i, \sigma)$.

O valor \hat{Y}_i será o valor predicado pola estimación do modelo, é dicir, $\hat{Y}_i = a + bx_i$. Nótese en particular que $\hat{Y} = \bar{Y} = a + b\bar{x}$.

Así, despois de multiplicar por n , a variabilidade é estimada mediante

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n (Y_i - \hat{Y}_i)^2.$$

Se agora denotamos

$$\begin{aligned}SS_Y &= \sum_{i=1}^n (Y_i - \bar{Y})^2, \\SS_R &= \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = b^2 \sum_{i=1}^n (x_i - \bar{x})^2, \\SS_E &= \sum_{i=1}^n (Y_i - \hat{Y}_i)^2,\end{aligned}$$

entón, a expresión anterior pode escribirse como

$$SS_Y = SS_R + SS_E$$

(variabilidade total) (variabilidade debida á regresión) (variabilidade non explicada)

As fórmulas anteriores refírense ás “sumas de cadrados”. Se en lugar diso queremos as varianzas, simplemente hai que dividir polo tamaño mostral n :

$$s_Y^2 = \frac{1}{n} SS_Y, \quad s_R^2 = \frac{1}{n} SS_R, \quad s_E^2 = \frac{1}{n} SS_E.$$

Estas cantidades son unha estimación das varianzas teóricas obtidas no apartado anterior. Por outra banda, o coeficiente de determinación estímase mediante r^2 , de xeito que temos

$$r^2 = \frac{s_R^2}{s_Y^2} = \frac{SS_R}{SS_Y}.$$

Así, a estimación do coeficiente de determinación r^2 interprétase como a proporción da variabilidade da variable aleatoria Y que é explicada por X mediante o modelo de regresión.

Esta técnica de análise da varianza utilízase para comprobar se unha liña recta mostra unha cantidade significativa de variabilidade observada de Y . Se o suposto é que a regresión é válida, entón o que terá que suceder é que a maior parte da variabilidade terá que ser explicada por SS_R , sendo a parte non explicada pequena.

Obsérvase agora a equivalencia das seguintes condicións

$$\beta = 0 \Leftrightarrow \rho = 0,$$

é dicir, toda a variabilidade é aleatoria (non explicada), e por tanto non hai regresión linear. Así pois, o test que temos que facer para comproba-la validez do modelo é

$$H_0: \rho = 0, \quad H_1: \rho \neq 0.$$

Baixo as hipóteses anteriores, este contraste emprega dous estatísticos que pasamos a describir a continuación. En primeiro lugar, para SS_Y hai n datos e un valor estimado, \bar{Y} , o que deixa $n - 1$ graos de liberdade.

- Para SS_E hai n datos, pero dous valores estimados, a e b , o que nos deixa $n - 2$ graos de liberdade. Así, empregamos como estatístico o *cadrado medio do erro*:

$$MS_E = \frac{SS_E}{n - 2}.$$

- Iso significa que para SS_R queda un só grao de liberdade. O estatístico empregado é pois o *cuadrado medio da regresión*: $MS_R = \frac{SS_R}{1}$.

No suposto de que a hipótese nula sexa certa, o estatístico

$$\frac{MS_R}{MS_E} = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2 / (n-2)} \sim F_{1, n-2}$$

segue unha distribución F de Snedecor con $(1, n-2)$ graos de liberdade.

Se a hipótese nula é certa, o valor observado no estatístico estará próximo a 1. Noutro caso será moito maior e rexeitarase a hipótese nula se o valor é demasiado grande. Trátase por tanto de facer un *contraste unilateral dereito*.

Os cálculos necesarios para empregar ANOVA á hora de contrastar $H_0: \rho = 0$ (non hai regresión linear), dispóñense nunha táboa como a seguinte:

variabilidade	g.l.	SS	MS	cociente
regresión	1	$SS_R = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = nr^2 s_Y^2$	$MS_R = \frac{SS_R}{1}$	$F_{1, n-2} = \frac{MS_R}{MS_E}$
erro	$n-2$	$SS_E = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = n(1-r^2) s_Y^2$	$MS_E = \frac{SS_E}{n-2}$	
total	$n-1$	$SS_Y = \sum_{i=1}^n (Y_i - \bar{Y})^2 = n s_Y^2$		

Cando, despois de face-lo anterior contraste de hipóteses, cheguemos á conclusión de que se rexeita a hipótese nula H_0 , iso quererá dicir que unha parte significativa da variabilidade de Y se pode explicar mediante o modelo de regresión linear. Iso non quere dicir que o modelo linear sexa o mellor para explicar dita variabilidade, senón que é razoable empregar-lo modelo para explicala.

Problema 6.4. Realízase un experimento para estuda-la relación entre a altura e a lonxitude da concha de *Patelloida pygmaea* (en mm). Téñense os seguinte datos:

altura	lonxitude	altura	lonxitude	altura	lonxitude	altura	lonxitude
0.9	3.1	1.9	5.0	2.1	5.6	2.3	5.8
1.5	3.6	1.9	5.3	2.1	5.7	2.3	6.2
1.6	4.3	1.9	5.7	2.1	5.8	2.3	6.3
1.7	4.7	2.0	4.4	2.2	5.2	2.3	6.4
1.7	5.5	2.0	5.2	2.2	5.3	2.4	6.4
1.8	5.7	2.0	5.3	2.2	5.6	2.4	6.3
1.8	5.2	2.1	5.4	2.2	5.8	2.7	6.3

Estima-la recta de regresión da lonxitude como función da altura. Calcula-lo coeficiente de determinación e interpreta-lo seu valor. ¿Hay evidencia estatística de que o modelo de regresión linear é válido?

Solución. Chamemos X á altura e Y á lonxitude. Organizámo-los cálculos nunha táboa.

X	Y	X^2	XY	Y^2	
0.90	3.10	0.81	2.79	9.61	
1.50	3.60	2.25	5.40	12.96	
1.60	4.30	2.56	6.88	18.49	
1.70	4.70	2.89	7.99	22.09	
1.70	5.50	2.89	9.35	30.25	
1.80	5.70	3.24	10.26	32.49	
1.80	5.20	3.24	9.36	27.04	
1.90	5.00	3.61	9.50	25.00	
1.90	5.30	3.61	10.07	28.09	
1.90	5.70	3.61	10.83	32.49	
2.00	4.40	4.00	8.80	19.36	
2.00	5.20	4.00	10.40	27.04	
2.00	5.30	4.00	10.60	28.09	
2.10	5.40	4.41	11.34	29.16	
2.10	5.60	4.41	11.76	31.36	
2.10	5.70	4.41	11.97	32.49	
2.10	5.80	4.41	12.18	33.64	
2.20	5.20	4.84	11.44	27.04	
2.20	5.30	4.84	11.66	28.09	
2.20	5.60	4.84	12.32	31.36	
2.20	5.80	4.84	12.76	33.64	
2.30	5.80	5.29	13.34	33.64	
2.30	6.20	5.29	14.26	38.44	
2.30	6.30	5.29	14.49	39.69	
2.30	6.40	5.29	14.72	40.96	
2.40	6.40	5.76	15.36	40.96	
2.40	6.30	5.76	15.12	39.69	
2.70	6.30	7.29	17.01	39.69	
Σ	56.60	151.10	117.68	311.96	832.85

Entón temos $n = 28$ datos e

$$\bar{X} = \frac{56,6}{28} = 2,021,$$

$$\bar{Y} = \frac{151,10}{28} = 5,396,$$

$$s_X^2 = \frac{117,68}{28} - 2,021^2 = 0,117,$$

$$s_Y^2 = \frac{832,85}{28} - 5,396^2 = 0,623,$$

$$s_{XY} = \frac{311,96}{28} - 2,021 \cdot 5,396 = 0,233.$$

Obtemos $b = 0,233/0,117 = 1,996$ e $a = 5,396 - 1,996 \cdot 2,020 = 1,361$ co que a ecuación da recta de regresión é

$$y - 5,396 = 1,996(x - 2,0214),$$

ou ben,

$$y = 1,361 + 1,996x.$$

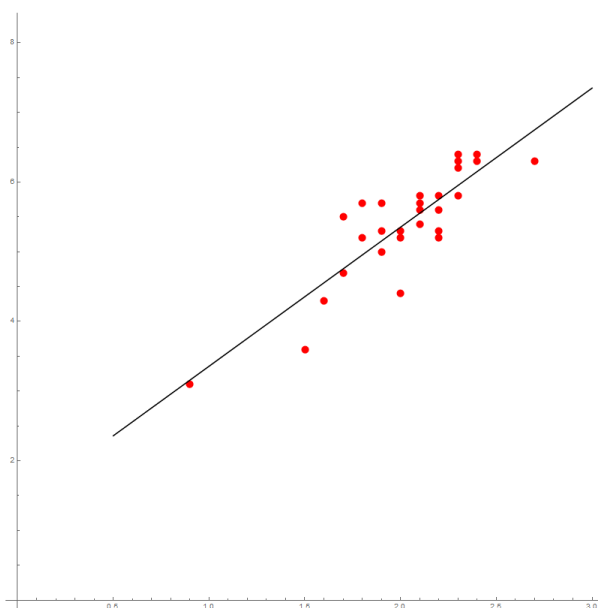


Figura 6.6: Os puntos e a súa recta de regresión

A estimación do coeficiente de correlación é

$$r = \frac{0,233}{\sqrt{0,117 \cdot 0,623}} = 0,8638,$$

de xeito que a calidade da aproximación parece moderada.

A estimación do coeficiente de determinación é $r^2 = 0,746$. Isto interprétase do seguinte xeito: o 74.6 % da variabilidade da variable Y está explicada polo modelo de regresión.

Para asegurarnos, intentaremos dar evidencia significativa de que o modelo de regresión é válido. Isto significa face-lo contraste de hipóteses

$$H_0: \rho = 0, \quad H_1: \rho \neq 0.$$

Empregamos pois a técnica de análise da varianza, ANOVA. Os datos necesarios están recollidos na seguinte táboa:

variabilidade	g.l.	SS	MS	cociente
regresión	1	$SS_R = 28 \cdot 0,864^2 \cdot 0,623 = 13,02$	$MS_R = 13,02$	76,42
erro	26	$SS_E = 28(1 - 0,864^2)0,623 = 4,43$	$MS_E = \frac{4,43}{26} = 0,17$	
total	27	$SS_Y = 28 \cdot 0,623 = 17,45$		

Como $P = P(F_{1,26} \geq 76,42) < 0,01$ é un número moi pequeno (de feito, empregando software estatístico temos $P = 3,2 \cdot 10^{-9}$), rexeitámo-la hipótese nula. Concluimos que hai evidencia significativa de que o modelo de regresión linear é válido. \square

6.2.2. Intervalos de estimación

Por completitude incluímos nesta sección a estimación por intervalos de diversos valores que apareceron no noso modelo de regresión linear.

En primeiro lugar pódese ver que unha estimación puntual da desviación típica do erro σ^2 vén dada por

$$\hat{\sigma}^2 = MS_E = \frac{SS_E}{n-2} = \frac{1}{n-2} \sum_{i=1}^n (Y_i - a - bx_i)^2,$$

tendo o estatístico

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2}$$

unha distribución χ^2 con $n-2$ graos de liberdade.

Tomemos un nivel de significación α .

- Ordenada na orixe:

$$\alpha = a \pm t_{n-2, \alpha/2} \frac{\sqrt{MS_E}}{S_X} \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}.$$

- Pendente da recta de regresión:

$$\beta = b \pm t_{n-2, \alpha/2} \frac{\sqrt{MS_E}}{S_X}.$$

- Resposta media para un valor de X dado:

$$\mu_{Y|X=x} = \hat{Y} \pm t_{n-2, \alpha/2} \sqrt{MS_E} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{s_X^2}}.$$

- Intervalo de predicción da resposta individual para un valor de X dado:

$$\hat{Y} \pm t_{n-2, \alpha/2} \sqrt{MS_E} \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{s_X^2}}.$$

Capítulo 7

Problemas para as clases interactivas

7.1. Intervalos de confianza

1. Nunha mostra de tamaño 30 mediuse a porcentaxe de aumento de alcohol en sangue tras beber catro cervexas. Obtívose $\bar{X} = 41,2$ (media) e $s = 2,1$ (cuasi-desviación típica).
 - a) Calcular un intervalo de confianza do 90 % para a porcentaxe media de aumento en tódalas persoas que beben catro cervexas;
 - b) Se se calcula un intervalo de confianza do 95 % para μ , ¿será máis ou menos amplo có anterior?

[Milton 6.3.7]

2. As granxas de patos contaminan a agua debido ó nitróxeno en forma de “ácido úrico”. A seguinte é unha mostra aleatoria de nove observacións da variable X , número de kilos de nitróxeno producidos por granxa e día.

4.9	5.8	5.9
6.5	5.5	5.0
5.6	6.0	5.7

Supoñendo que X é normal, construír un intervalo de confianza do 99 % para a media poboacional μ .

[Milton 6.3.9]

3. A calor parece afecta-la mobilidade dos caracois. En 20 caracois sometidos a unha temperatura de 29°C observamos unha distancia media percorrida de $\bar{X} = 4,855\text{cm}$, con $s_{n-1} = 0,7178$. Dar un intervalo de confianza ($\alpha = 5\%$) para a distancia media percorrida por un caracol.

[Milton 6.3.6 p. 222]

4. Estas son as alturas (en metros) de vinte piñeiros da especie “Pinus strobus”. Estima-la media desa especie de piñeiros cun nivel de confianza do 95 %.

17.16	22.00	10.08	15.00
7.02	10.67	11.16	10.92
11.10	4.05	15.93	7.22
8.19	16.45	7.38	10.00
14.10	10.26	11.96	10.00

[Milton 6.3.1]

5. Queremos estima-lo peso medio ó nacer (en Kg) de fillos de mulleres adictas á heroína. Nun estudio previo obtívose que $\sigma = 2,5$. Queremos deseña-lo experimento de modo que o nivel de confianza sexa do 95 %, e que o erro de estimación non supere 1Kg. ¿Que tamaño de mostra necesitamos?

[Milton 6.6.1 p. 236]

6. No río Mississippi estudouse en 61 lugares a variable X , anchura de terreno inundable, obténdose $\bar{X} = 3400$ metros e $s_{n-1} = 100$ metros. Dar un intervalo de estimación para a desviación típica de X cun nivel de confianza do 90 %.

[Milton 7.1.7 p. 253]

7. Nun recuento no microscopio contabilizáronse 200 leucocitos, dos cales 125 eran neutrófilos. Dar un intervalo de confianza do 90 % para a proporción de neutrófilos en sangue.

[Milton 8.2.3 p. 266]

8. Nun estudo sobre obesidade infantil averíguase que a idade media de inicio da enfermidade dunha mostra de 26 nenos é de 4 anos, cunha desviación típica mostral de 1.5 anos. Determinar un intervalo de confianza do 95 % para a desviación típica da poboación.

[Milton Exemplo 7.1.6]

9. ¿Que tamaño de mostra faría falla para estima-la proporción de mortes debidas a un problema cardíaco, se traballamos cun nivel de significación do 5 %, e non queremos que o erro de estimación supere o 2 %?

[Milton 8.3.2 p. 270]

10. Un investigador médico quere estima-lo nivel medio de colesterol en homes de idade avanzada. A estimación debe ter unha precisión de 6mg/dl ou menos, cun 95 % de confianza. Ademais, o investigador cre, por estudos previos, que a desviación típica do colesterol na poboación ronda os 40mg/dl. ¿Que tamaño de mostra debe tomar?

[Samuels 6.4.2]

11. Nun estudo atopouse que 40 de 400 estudantes eran zurdos. Construír un intervalo de confianza do 90 % para a proporción de estudantes zurdos na poboación.

[Samuels 9.3.1]

12. Unha bodega produce 720000 botellas de viño cada ano e desexa estima-la proporción de botellas que teñen o corcho defectuoso (o viño estropéase se hai un fallo no corcho). Nun estudo previo calcúlase que esta proporción ronda o 4 %, pero agora queremos, cun nivel de confianza do 90 %, que o erro de estimación non supere o 1 %. ¿Cantas botellas de viño debemos comprobar?

[Samuels 9.S.6]

7.2. Contrastes de hipóteses

1. Sospéitase que o insecticida DDT provoca diminución no grosor das cáscaras dos ovos dos paxaros. Para combrobar isto, alimentouse a 16 gabiáns cunha mistura que contiña 15ppm de DDT, e atopouse unha diminución do grosor do 8 %. A desviación típica mostral foi de $s = 0,05$. Contrasta-la hipótese de que houbo unha diminución no grosor en toda a poboación (nivel de confianza do 95 %).

[Milton 6.5.4 p. 233]

2. Realizouse un experimento para estuda-lo efecto do exercicio físico no nivel de colesterol de pacientes obesos. En 80 pacientes sometidos a un réxime específico de actividade, observouse unha diminución media do nivel de colesterol de $\bar{X} = 27$ puntos. A desviación estándar foi de $s = 18$. ¿Pode afirmarse, cun nivel de confianza do 90 %, que ese réxime provoca, en media, unha diminución superior a 25 puntos?

[Milton 6.5.7 p. 234]

3. A concentración media de dióxido de carbono no aire é do 0.035 %. Preténdese demostrar que inmediatamente por riba da superficie do chan dita concentración é maior. Analizáronse 144 mostras de aire seleccionado aleatoriamente e tomadas á distancia de 30cm do chan. Resultou unha media mostral do 0.09 % e unha cuasi-desviación típica mostral do 0.25 %. ¿Cal é o valor P do contraste? ¿Comprobouse estatisticamente o argumento establecido?

[Milton 6.5.5 p. 233]

4. En certa especie de vagalumes, a luz que producen consta dun escintileo curto seguido dun período de repouso. Quérese probar que o período de repouso ten unha duración media de menos de catro segundos. Nunha mostra de 16 insectos obtivemos unha media de 3.77 segundos, con $s = 0,30$ segundos. Por outro lado, dámonos conta de que un erro de tipo I non ten consecuencias fatais, así que fixamos un $\alpha = 10\%$ bastante alto. ¿Apoian os datos experimentais a nosa suposición sobre o escintileo?

[Milton 6.5.7 p. 232]

5. Ó estuda-lo crecemento de abetos, sábese que a varianza poboacional acostuma ser $1,56\text{cm}^2$. Non obstante, en 50 árbores crecidos en condicións de seca observamos unha cuasi-desviación típica de 0,375cm. ¿Afectou a seca ó parámetro σ ? Dar un intervalo de confianza do 95 % para a desviación típica da poboación.

[Milton 7.2.4 p. 256]

6. A concentración sanguínea de calcio nos mamíferos acostuma ser 6mg/100ml. A desviación típica debe ser de 1mg/100ml, xa que unha variabilidade maior ocasiona trastornos de coagulación. Nunha serie de nove probas realizadas a un paciente, atopouse unha concentración media de 6.2 e unha cuasi-desviación típica de 2. Tomando un nivel de significación $\alpha = 0,05$, ¿hai evidencia de que a desviación típica sexa maior da normal?

[Milton 7.2.2 p. 256]

7. Estímase xeralmente que o 90 % dos enfermos de cancro de pulmón morren no prazo de 3 anos. Nun estudo recente no que se proban uns novos tratamentos, atopouse que 128 pacientes morreron dun total de 150 enfermos. ¿Pode dicirse que hai probas suficientes de que o emprego dos novos métodos de tratamento reduciron a taxa de falecementos?

[Milton 8.4.4 p. 273]

8. Un 20 % dos enfermos de corazón tratados cronicamente con digoxina sofre unha reacción adversa. Para evitalo, a 30 pacientes asocióuselles outro medicamento, e conseguíuse que só tres tivesen a reacción. ¿Pode afirmarse que o tratamento é eficaz cun nivel de confianza do 99 %?

[Milton 9.7]

9. O método habitual para trata-la leucemia mieloblástica aguda consiste en somete-lo paciente a quimioterapia intensiva no momento do diagnóstico. Historicamente, isto produciu unha taxa de remisión do 70 %. Estudando un novo método de tratamento utilizáronse 50 voluntarios. ¿Cantos dos pacientes deberían ter remitido para que os investigadores puidesen afirmar, con nivel de significación $\alpha = 0,025$, que o novo método produce remisións máis altas có antigo?

10. Os votos en contra da construción dunha presa nunha mostra de 500 persoas foi de 270. Estima-la proporción de persoas que están en contra en toda a poboación, cun nivel de confianza do 95 %.

[Milton Exemplo 8.4.1]

11. Estase probando a eficacia dun tipo de exercicio para mellora-los síntomas da artrite reumatoide. O grupo no que se proba dito tratamento é de 160 pacientes. Para un nivel de significación do 2,5 %, ¿cantos pacientes terían que mellorar para que se poida afirmar que a porcentaxe de pacientes que melloran é superior ó 50 %?

7.3. Contrastes de hipóteses para dúas poboacións

1. Comprobase o peso de ovos de tartaruga en dúas illas diferentes. Suponse que a variable é normal. Á vista dos datos obtidos en dúas mostras aleatorias, ¿hai evidencia de que os ovos na illa “Malabar” son máis pesados cós da illa “Grande-Terre” cun nivel de significación do 1 %?

Datos da illa “Grande-Terre”: Tamaño da mostra $n_1 = 31$; peso medio $\bar{X}_1 = 64,0g$; cuasi-desviación típica $s_1 = 6,5g$.

Datos da illa “Malabar”: Tamaño da mostra $n_2 = 148$; peso medio $\bar{X}_2 = 82,7\text{g}$; cuasi-desviación típica $s_2 = 3,6\text{g}$.

(Facer un contraste de hipóteses para a igualdade das varianzas para poder determinar se podemos asumir que ambas sexan iguais.)

[Milton 9.4.3 p. 311]

2. Ó estuda-la velocidade de voo de dúas especies de paxaros, obtivémo-los seguintes datos:

- (*Haematopus palliatus*): $n_1 = 9$, $\bar{X}_1 = 26,05$, $s_1 = 6,34$;
- (*Pelecanus occidentalis*): $n_2 = 12$, $\bar{X}_2 = 30,19$, $s_2 = 3,20$;

Face-lo contraste necesario para saber se as varianzas poboacionais se poden supoñer iguais. ¿Hai evidencia de que a velocidade de voo das dúas especies de paxaros sexa diferente? (Para todo o problema, tomar un nivel de confianza do 95 %.)

[Milton 9.2.1 p. 298]

3. Estudouse nunha mostra de $n_1 = 33$ homes novos fumadores a idade media á que empezan a fumar, obténdose $\bar{X}_1 = 11,3$ anos. A cuasi-varianza mostral foi de 4 anos. O mesmo estudo en mozas deu lugar ós seguintes datos: $n_2 = 14$, $\bar{X}_2 = 12,6$, $s_2^2 = 3,5$. Pídese, cun nivel de significación $\alpha = 5\%$:

- a) Facer unha proba F para concluír que podemos supoñer $\sigma_1^2 = \sigma_2^2$;
- b) Dar un intervalo de estimación para a diferenza de medias poboacionais entre mozos e mozas.

[Milton 9.3.11 p. 309]

4. Un laboratorio quere compara-los efectos secundarios dun medicamento novo cos do produto da competencia. Usaremos un nivel de significación do 1%. Obtivéronse os seguintes datos sobre a porcentaxe de persoas que presentaban diarrea:

	Laboratorio	Competencia
Número de suxeitos	465	195
Número de casos de diarrea	9	1

- a) ¿Podemos afirmar que as porcentaxes son significativamente diferentes?
- b) Dar un intervalo de confianza para a diferenza de porcentaxes.

[Milton 8.6.6 p. 285]

5. En 1970 fixéronse 759 análises de sangue e atopáronse 46 casos de infección. En 1975 outro estudo semellante descubriu 109 infeccións en 838 análises. Baseándose nestas dúas mostras, ¿podemos estar seguros de que a proporción de casos de infección aumentou en máis de 6 puntos porcentuais nesos cinco anos? (Usar nivel de confianza do 90 %.)

[Milton 8.6.4]

6. A partir dos corenta anos, o cancro de mama pode detectarse a través dunha mamografía. Comprobamos que en 31 mulleres novas afectadas (idade 40-49 anos) houbo 6 casos descubertos a través de mamografía. Por outra parte, nun grupo de 101 mulleres de máis idade, a mamografía foi eficaz en 38 casos. Cun nivel de confianza do 95 %, ¿podemos afirmar que a mamografía é menos eficaz nas mulleres novas?

[Milton 8.6.3 p. 285]

7. Para ver se un medidor portátil de glucosa é útil para os diabéticos, mediuse para cada paciente o nivel de glucosa en sangue antes de aprender a usalo, e unhas semanas despois. Nunha mostra aleatoria de 36 individuos atopouse unha diferenza de 2.78mmol/l entre “antes” e “despois”, con cuasi-desviación típica das diferenzas igual a 6.05. ¿Quere dicir isto que o medidor é efectivo para axudar a reduci-los niveis de glucosa?

[Milton 9.5.3 p. 319]

8. Os datos de temperatura en 1000 estacións meteorolóxicas en todo o mundo deron unha temperatura media de 57 graos Fahrenheit en 1950, e de 57.6 en 1988, con $s_D = 4,1$. ¿Quere isto dicir que a temperatura media do globo aumentou? (Empregalo valor P .) Dar un intervalo para o aumento global medio (para un nivel de confianza do 90 %).

[Milton 9.5.5 p. 320]

7.4. Problemas de repaso de estimación e contraste de hipóteses

1. Para que un peixe sobreviva, a cantidade de osíxeno disolto na auga non debe ter unha desviación típica maior cá 1.2 partes por millón. Tomamos mostras de auga en 25 lugares aleatoriamente escollidos dun lago e obtemos $s = 1,7$ ppm. ¿Evidencia isto que a variabilidade do osíxeno aumentou por riba do parámetro aceptable $\sigma = 1,2$?

[Milton 7.2.3 p. 256]

2. Nun estudo sobre rexeneración de células nerviosas en monos *rhesus* mediuse o contido en creatinina fosfato na parte esquerda e na parte dereita da espiña dorsal (medido en mg de CF por cada 100g de tecido). Para un nivel de significación do 10 %, ¿existe evidencia significativa dunha diferenza na cantidade de CF entre os dous datos? Os datos son os seguintes:

Animal	1	2	3	4	5	6	7	8
Lado dereito	16.3	4.8	10.9	14.2	16.3	9.9	29.2	22.4
Lado esquerdo	11.5	3.6	12.5	6.3	15.2	8.1	16.6	13.1

[Samuels p. 333]

3. Crese que a maioría dos fumadores empezan a fumar despois dos 18 anos. Nunha mostraxe con 60 individuos, atopouse que o 49 % empezou a fumar despois desa idade.

- a) Decidir se hai evidencia de que na poboación a proporción de fumadores que empeza despois dos 18 é menor có 50 % (cun nivel de significación do 1 %).
- b) Explica-las consecuencias económicas e sanitarias de cometer un erro de tipo I ou un erro de tipo II.

[Milton 6.4.6 p. 226]

4. Existe a teoría de que a vitamina C é beneficiosa no tratamento do cancro. Os que a defenden din que hai unha melloría superior ó 4% de casos. Fixemos dous grupos independentes de 75 individuos cada un. Ós primeiros démoslle 10g diarios de vitamina C; ós outros, nada. Ó cabo de catro semanas, no primeiro grupo 47 pacientes presentaron algunha melloría, mentres que este número foi soamente de 43 no segundo grupo. Pídese face-lo contraste $H_0: p_1 - p_2 \leq 0,04$ e interpreta-lo resultado (emprega-lo valor P).

[Milton 8.6.1 p. 280 e 8.6.2 p. 281]

5. Estase probando a eficacia de dous tipos de exercicio para mellora-los síntomas da artrite reumatoide. O primeiro tratamento (T1) foi probado en 150 pacientes con esta enfermidade obtendendo que 87 deles melloran tras un mes de práctica. O segundo tratamento (T2) foi probado en 160 pacientes dos que 72 melloraron tras un mes de práctica. ¿Podemos asegurar que hai evidencia significativa de que a proporción de pacientes que melloran co tratamento T1 é superior á do T2? Realiza-lo correspondente contraste de hipóteses.

7.5. Probas de homoxeneidade e independencia

1. Investígase a eficacia dunha nova vacina contra a gripe. Elíxese unha mostra de 900 persoas, e clasifícanse segundo que foran ou non vacinadas, e segundo contraeran a gripe durante o último ano ou non. Pídese, cun nivel de confianza do 95 %, decidir se hai asociación ou non entre as dúas variables.

Vacinado \ gripe	si	non
si	150	200
non	300	250

[Milton 12.1.2 p. 449]

2. Cremos que existe relación entre o número de cloroplastos das follas das árbores e o nivel de SO_2 no aire. Selecciónanse 60 árbores, e clasifícanse en función do nivel de dióxido de azufre da súa zona e o nivel de cloroplastos das súas follas. Obtéñense os seguintes datos:

SO_2 \ Cloroplastos	alto	normal	baixo
alto	5	4	13
normal	5	10	5
baixo	7	9	2

- a) ¿Trátase dunha proba de independencia ou de homoxeneidade?
 b) ¿Que conclusións poden sacarse dos datos? Enuncia a hipótese nula apropiada e razoa en función do valor P obtido.

[Milton 12.2.6 p. 457]

3. Co obxectivo de provoca-la unión dos ósos en fracturas, aplícanse campos electromagnéticos pulsantes. Nunha mostra de 62 fracturas de tibia, 26 de húmero, e 18 de fémur, observouse que o tratamento só tivo éxito en 34, 16, e 10 delas, respectivamente.
- Construí-la táboa de continxencia axeitada.
 - Á vista dos resultados obtidos na mostra, ¿pódese concluír que o éxito do tratamento depende do tipo de óso que se está tratando?
4. Realízase un pequeno estudo piloto para determinar se hai asociación entre a aparición de leucemia e os antecedentes de alerxia. Selecciónase unha mostra de 19 pacientes con leucemia e outro grupo de control de 17 persoas, e determínase se hai antecedentes de alerxia ou non.

grupo \ antecedentes	si	non
paciente	17	2
control	5	12

Calcula-la frecuencia esperada para cada celda e contrastar se a distribución de casos de alerxia é homoxénea nos dous grupos. Explica-la resposta baseándose no valor P do contraste.

[Milton 12.1.5 p. 449]

5. Nun estudo sobre quimioterapia no cancro de pulmón administráronse simultaneamente catro medicamentos a 16 pacientes, mentres que a outro grupo de 11 pacientes déronse os medicamentos de xeito secuencial. Observouse unha resposta positiva ó tratamento en 11 pacientes do primeiro grupo, e en 3 dos tratados secuencialmente. ¿Proporcionan estes datos evidencia de que unha forma de tratamento é superior á outra?

[Samuels 10.2.10]

7.6. Regresión linear e ANOVA

1. Realízase un estudo para estima-la relación entre o índice de obesidade X e a taxa metabólica en repouso Y . A partir dos datos de 43 individuos obtemos

$$\begin{aligned} \sum X &= 1482,5; & \sum Y &= 10719; \\ \sum X^2 &= 53515,25; & \sum Y^2 &= 2736063; & \sum XY &= 379207,5. \end{aligned}$$

- a) ¿Que taxa metabólica correspondería a un índice de obesidade $X = 40$?

- b) Calcular e interpreta-lo coeficiente de determinación.
- c) Contrasta-lo modelo de regresión linear.

[Milton 11.3.4 p. 414]

2. A seguinte táboa recolle os datos de presións sistólicas (P) de cinco individuos en función da súa idade (t):

t idade (anos)	20	30	40	50	60
P presión (mm Hg)	125	128	131	133	138

- a) ¿Que ecuación linear nos permite estimar P para un individuo de 25 anos?
 - b) Calcula-lo coeficiente de determinación e interpreta-lo resultado.
 - c) Contrasta-lo modelo de regresión linear.
3. Realizouse un experimento para estima-la concentración plasmática Y dunha substancia a partir da súa concentración X na saliva. Os datos experimentais foron:

X	7.4	7.5	8.5	9.0	11.0	13.0	14.0	14.5	16.0	17.0
Y	30.0	25.0	31.5	27.5	40.2	48.0	52.0	54.0	56.5	58.0

Calcula-la recta de regresión e contrasta-lo modelo de regresión linear (ANOVA).

4. A cantidade de arsénico no arroz (variable Y , en $\mu g/kg$) parece estar relacionada coa de silicio na palla de arroz (variable X , en g/kg). Ó estudar 32 plantas obtémo-os seguintes datos:

$$\bar{X} = 29,85, s_X = 10,04, \bar{Y} = 122,25 s_Y = 44,50, r = -0,556.$$

- a) ¿Que cantidade de arsénico estimamos cando $X = 12$?
- b) Calcula-la varianza residual dos erros de estimación.
- c) ¿Que proporción de varianza da concentración de arsénico está explicada pola relación linear co contido de silicio?

[Samuels p. 505]

5. Aplicáronse dous cuestionarios a 670 persoas: un medía o nivel de estrés ó que estiveran sometidas X , e o outro detectaba posibles trastornos de saúde Y . Ó calcula-lo coeficiente de correlación de Pearson obtívose $r = 0,24$. ¿É compatible este resultado coa hipótese $\rho = 0$? (tomar $\alpha = 5\%$)
6. Déronse distintas doses dunha substancia velenosa a sete grupos de 26 ratos, e observáronse os seguintes resultados:

X doses (mg)	4	6	8	10	12	14	16
Y número de mortes	1	3	6	8	14	16	20

- a) Calcula-la ecuación da recta de mínimos cadrados axustada a estes datos.

- b) Estima-lo número de mortes nun grupo de 26 ratos que recibiron unha dose de 7mg deste veneno.
- c) Contrasta-lo modelo de regresión linear.
7. Lévese a cabo un estudo sobre as características corporais e o modo de actuar de levantadores de peso olímpicos. Estúdanse as variables X , peso corporal, e Y , mellor levantamento, obtendo:

X	134	138	154	178	176	190	190	205	205	206
Y	185	238	260	290	312	336	339	341	358	359

- a) Debuxa-la nube de puntos. Baseándose nela, ¿pódese esperar que b sexa positivo ou negativo?
- b) Calcular e interpreta-lo coeficiente de determinación.
- c) Comproba-la idoneidade do modelo de regresión linear. Se é axeitado calcula-la liña de regresión de X sobre Y , estima-lo mellor levantamento dun atleta que pesa 200 libras.

[Milton 11.4.1]

Capítulo 8

Exames resoltos

8.1. Exame 1

Problema 8.1. Ó estuda-la coagulación do sangue utilízase a variable normal X , tempo parcial activado en segundos da tromboplastina. Os valores seguintes representan unha mostra aleatoria de 10 observacións sobre X para un determinado paciente:

45 40 47 46 42 50 47 48 49 49.

1. Construír un intervalo para o tempo parcial medio da tromboplastina para ese paciente, cun nivel de confianza do 99 %.
2. Se a varianza poboacional é 9, ¿cal ten que se-lo tamaño da mostra para que a diferenza entre a media mostral e a media poboacional sexa como moito de ± 1 segundo, cun nivel de confianza do 99 %?

Solución. Sexa pois X “tempo parcial activado en segundos da tromboplastina”.

Para o primeiro apartado temos que calcular un intervalo de confianza para a media (Subsección 2.2.3) empregando o estatístico

$$\frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}},$$

que segue unha distribución t_{n-1} . Despexando μ da inecuación

$$\left| \frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}} \right| \leq t_{n-1, \alpha/2}$$

obtense a fórmula

$$\bar{X} \pm t_{n-1, \alpha/2} \frac{s_{n-1}}{\sqrt{n}}.$$

Temos $n = 10$. Organizámo-los cálculos para calcula-la media e cuasi-varianza mostral:

X	X^2
45	2025
40	1600
47	2209
46	2116
42	1764
50	2500
47	2209
48	2304
49	2401
49	2401
Σ	463 21529

De aquí obtemos $\bar{X} = 463/10 = 46,3$, $s_n^2 = 21529/10 - 46,3^2 = 9,21$, e así, $s_{n-1} = \sqrt{\frac{10}{9} 9,21} = 3,20$.

Nivel de significación $\alpha = 0,01$. Buscámo-lo valor $t_{9,0,005} = 3,25$ nas táboas. Aplicando a fórmula

$$46,3 \pm 3,25 \frac{3,20}{\sqrt{10}} = 46,3 \pm 3,29,$$

de onde se deduce o intervalo $[43,01, 49,59]$.

Conclusión: cun nivel de confianza do 99%, o tempo parcial activado medio da tromboplastina atópase entre 43.01 e 49.59 segundos.

Para o segundo apartado témo-lo dato $\sigma^2 = 9$. Como a varianza poboacional é coñecida (Subsección 2.2.2), empregámo-lo estatístico

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}},$$

que ten distribución normal estándar. Despexando μ da inecuación

$$\left| \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right| \leq Z_{\alpha/2}$$

obtémo-la fórmula

$$\bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

A estimación do erro (Observación 2.5) é $Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$, e queremos $Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \epsilon$, onde ϵ é o valor fixado polo problema. Despexando n obtense $n \geq (Z_{\alpha/2} \sigma / \epsilon)^2$.

O nivel de confianza é $\alpha = 0,01$. Mirando as táboas obtemos $Z_{0,005} = 2,58$. Neste caso $\epsilon = 1$. Substituíndo na fórmula $n \geq (2,58 \cdot 3/1)^2 = 59,7$.

Conclusión: para que a diferenza entre a media mostral e a media poboacional no tempo parcial activado en segundos da tromboplastina sexa como moito de ± 1 segundo cun nivel de confianza do 99%, teriamos que tomar unha mostra de polo menos 60 elementos. \square

Problema 8.2. Estase a probar un antibiótico chamado DOXICICLINA para previr a “diarrea do viaxeiro”. O fármaco foi probado sobre 64 voluntarios que foron a Kenya. A unha metade déuselle doxiciclina e á outra un placebo. Dos que recibiron doxiciclina, 24 libráronse do trastorno, mentres que só 16 dos do outro grupo se libraron.

1. Construír un intervalo de confianza do 95 % para a diferenza entre as porcentaxes de protección entre aqueles que utilizaron doxiciclina e os que non a utilizaron. Interpreta-lo intervalo.
2. ¿Pódese asegurar que a doxiciclina contribúe a proporcionar protección contra a diarrea do viaxeiro? Explicalo sobre a base do valor P.

Solución. As variables aleatorias a considerar son X , non ter diarrea do viaxeiro entre voluntarios que tomaron doxiciclina, e Y , non ter diarrea do viaxeiro entre voluntarios que tomaron placebo.

Para o primeiro apartado temos que calcular un intervalo de confianza para a diferenza de porcentaxes (Sección 4.3) empregando o estatístico

$$\frac{(\widehat{p}_1 - \widehat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n_1} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{n_2}}},$$

que segue unha distribución normal estándar. Nótese que as poboacións non están emparelladas. O intervalo de confianza pedido obtense despregando $p_1 - p_2$ da desigualdade

$$\left| \frac{(\widehat{p}_1 - \widehat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n_1} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{n_2}}} \right| \leq Z_{\alpha/2},$$

de onde se obtén a fórmula

$$(\widehat{p}_1 - \widehat{p}_2) \pm Z_{\alpha/2} \sqrt{\frac{\widehat{p}_1(1-\widehat{p}_1)}{n_1} + \frac{\widehat{p}_2(1-\widehat{p}_2)}{n_2}}.$$

Temos como datos $n_1 = 32$, $\widehat{p}_1 = 24/32 = 0,75$, $n_2 = 32$, $\widehat{p}_2 = 16/32 = 0,5$.

Nivel de significación $\alpha = 0,05$. Buscamos na táboa $Z_{0,025} = 1,96$. Substituíndo na fórmula:

$$(0,75 - 0,5) \pm 1,96 \sqrt{\frac{0,75(1-0,75)}{32} + \frac{0,5(1-0,5)}{32}} = 0,25 \pm 0,229,$$

de onde se obtén o intervalo $[0,021, 0,479]$.

Conclusión: cun nivel de confianza do 95 %, a diferenza de proporción de viaxantes a Kenya que non tiveron a diarrea do viaxeiro entre os que tomaron doxiciclina e os que tomaron placebo sitúase entre o 2.1 % e o 47.9 %.

Para a segunda cuestión temos que face-lo contraste de hipóteses

$$H_0: p_1 \leq p_2, \quad H_1: p_1 > p_2.$$

Este contraste ten cero como valor nulo (Sección 4.3). En consecuencia, agora temos que emprega-lo estatístico

$$\frac{\widehat{p}_1 - \widehat{p}_2}{\sqrt{\widehat{p}(1-\widehat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}},$$

que tamén segue unha distribución normal estándar, e onde

$$\widehat{p} = \frac{n_1\widehat{p}_1 + n_2\widehat{p}_2}{n_1 + n_2}.$$

Substituíndo temos, en primeiro lugar

$$\hat{p} = \frac{32 \cdot 0,75 + 32 \cdot 0,5}{32 + 32} = 0,625,$$

o cal nos dá o valor no estatístico

$$\frac{0,75 - 0,5}{\sqrt{0,625(1 - 0,625) \left(\frac{1}{32} + \frac{1}{32}\right)}} = 2,07.$$

Calculamos agora o valor P (Subsección 3.1.3) mirando a táboa da distribución normal: $P = P(z > 2,07) = 0,01923$. Temos que $1\% < P < 2,5\%$.

Conclusión: rexeitámo-la hipótese nula e concluímos que existe evidencia significativa, polo menos do 97.5%, de que a doxiciclina aumenta a proporción de viaxantes a Kenya que non teñen diarrea do viaxeiro fronte a aqueles que tomaron placebo. Por tanto, a doxiciclina contribúe a proporcionar protección contra a diarrea do viaxeiro. \square

Problema 8.3. A seguinte táboa representa as presións sanguíneas sistólicas (mm Hg) de 10 individuos alcohólicos rehabilitados, antes e despois de deixa-la bebida

Individuo	1	2	3	4	5	6	7	8	9	10
Antes	140	165	160	160	175	190	170	175	155	160
Despois	145	150	150	155	170	175	160	165	145	170

Supoñendo que as poboacións están distribuídas normalmente,

1. Estimar mediante un intervalo de confianza do 95% o cambio da presión sistólica que produce o abandono do alcohol. Interpretar o devandito intervalo.
2. ¿Hai evidencias suficientes, cun nivel de significación do 5%, para dicir que a presión sanguínea sistólica diminúe despois de deixa-la bebida?

Solución. As variables aleatorias a considerar son X , presión sanguínea sistólica dun alcohólico antes de deixa-la bebida, e Y , presión sanguínea sistólica dun alcohólico despois de deixa-la bebida. Obviamente trátase dun problema de comparación de dúas poboacións con mostras emparelladas (Sección 4.4), así que debemos toma-la variable diferencia $D = X - Y$.

O estatístico que temos que tomar é $\frac{\bar{D} - \mu_D}{s_D/\sqrt{n}}$, que segue unha distribución t_{n-1} . O primeiro que facemos é dispoñer os datos para calcula-los elementos da fórmula:

X	Y	D	D^2	
140	145	-5	25	
165	150	15	225	
160	150	10	100	
160	155	5	25	
175	170	5	25	
190	175	15	225	
170	160	10	100	
175	165	10	100	
155	145	10	100	
160	170	-10	100	
Σ	1650	1585	65	1025

Temos $n = 10$. Por tanto, $\bar{D} = 65/10 = 6,5$, $s_{n,D}^2 = 1025/10 - 6,5^2 = 60,25$, e $s_{n-1,D} = \sqrt{\frac{10}{9} 60,25} = 8,18$.

Como no primeiro apartado temos que calcular un intervalo de confianza, despxamos μ_D da desigualdade

$$\left| \frac{\bar{D} - \mu_D}{s_D/\sqrt{n}} \right| \leq t_{n-1, \alpha/2},$$

de onde obtemos $\bar{D} \pm t_{n-1, \alpha/2} \frac{s_D}{\sqrt{n}}$.

Nivel de significación $\alpha = 0,05$. Mirámo-lo valor $t_{9, 0,025} = 2,2622$ nas táboas. Substituíndo na fórmula anterior obtemos

$$6,5 \pm 2,26 \frac{8,18}{\sqrt{10}} = 6,5 \pm 5,85,$$

o que nos dá un intervalo $[0,64, 12,35]$.

Conclusión: cun nivel de confianza do 95 %, a diferenca media das presións sanguíneas sistólicas dun alcohólico rehabilitado entre antes e despois de deixa-la bebida sitúase entre 0.6 e 12.3mm Hg.

Para a segunda parte do exercicio, temos que face-lo seguinte contraste de hipóteses:

$$H_0: \mu_D \leq 0, \quad H_1: \mu_D > 0.$$

Como xa calculámo-los datos, substituímos no estatístico

$$\frac{6,5 - 0}{8,18/\sqrt{10}} = 2,51.$$

Pero agora necesitamos mirar na táboa $t_{9, 0,05} = 1,83$, que é menor ca 2.51.

Conclusión: rexeitamos H_0 e concluímos que hai evidencia significativa, ó 95 % de confianza, de que a presión sanguínea sistólica dun alcohólico rehabilitado diminúe despois de deixa-la bebida. \square

Problema 8.4. Diseñouse un estudo para analiza-la posible relación entre o medio no que viven e a incidencia de trastorno depresivo das persoas no paro. Seleccionáronse suxeitos pertencentes a medios rurais, semiurbanos e urbanos. De cada medio seleccionouse unha mostra aleatoria de 100 suxeitos no paro, obtendo que 12 do rural, 16 do semiurbano e 32 do urbano presentaban trastorno depresivo.

1. Construí-la táboa de continxencia axeitada. ¿Trátase dunha proba de independencia ou de homoxeneidade?
2. ¿Pode afirmarse, cun 1 % de nivel de significación, que na poboación de desempregados existe relación entre o tipo de medio no que se vive e padecer ou non trastorno depresivo?

Solución. Temos tres poboacións dependendo do medio no que viven, e a variable aleatoria Y ="incidencia de trastorno depresivo". En primeiro lugar construímo-la táboa de continxencia:

medio \ trastorno	si	non	tamaño
rural	12	88	100
semiurbano	16	84	100
urbano	32	68	100
Σ	60	240	300

O tamaño da mostra en cada medio está fixado polo investigador, trátase dunha proba de homoxeneidade para datos categóricos (Sección 5.2). Por tanto, temos que face-lo contraste de hipóteses:

$$H_0: p_{11} = p_{21} = p_{31}, p_{12} = p_{22} = p_{32}.$$

A continuación calculámo-los valores esperados no suposto de que houbose homoxeneidade nas poboacións mediante a fórmula $\widehat{E}_{ij} = \frac{n_i n_{\cdot j}}{n}$ (en verde), e tamén os valores $(n_{ij} - \widehat{E}_{ij})^2 / \widehat{E}_{ij}$ (en vermello), obtendo:

medio \ trastorno	si	non	Σ
rural	12 20 3.2	88 80 0.8	100
semiurbano	16 20 0.8	84 80 0.2	100
urbano	32 20 7.2	68 80 1.8	100
Σ	60	240	300

Finalmente aprovéitanse todas estas contas para calcula-lo valor no estatístico, (que consiste en suma-los valores vermellos), para obter 14.

O estatístico segue unha distribución χ^2 con $(3 - 1)(2 - 1) = 2$ graos de liberdade. Damos un nivel de significación $\alpha = 0,01$, así que índonos ás táboas obtemos $\chi_{2,0,01}^2 = 9,21$, que é menor ca 14.

Conclusión: rexeitámo-la hipótese nula, e concluimos que hai evidencia significativa, cun nivel de confianza do 99%, de que a incidencia de trastorno depresivo nas persoas en paro é distinto dependendo de se o medio no que viven é rural, semiurbano ou urbano. \square

Problema 8.5. Os seguintes datos corresponden a idade (X en anos) e a conduta agresiva (Y medida nunha escala de 0 a 10) dun grupo de 10 nenos, de entre 6 e 9 anos, elexidos ó azar

$$\sum X = 75, \sum Y = 49, \sum X^2 = 570,72, \sum Y^2 = 313, \sum XY = 345,2.$$

1. Estima-la recta de regresión que permita predicir o valor da conduta agresiva en función da idade do neno.
2. Calcula-lo coeficiente de determinación r^2 e interpreta-lo seu resultado.
3. Contrasta-lo modelo de regresión lineal.

Solución. Estamos chamando X á idade en anos, e Y á conducta agresiva dos nenos. Temos que calcula-la recta de regresión (Subsección 6.1.1) de Y sobre X .

Entón temos $n = 10$ datos e

$$\begin{aligned}\bar{X} &= \frac{75}{10} = 7,5, \\ \bar{Y} &= \frac{49}{10} = 4,9, \\ s_X^2 &= \frac{570,72}{10} - 7,5^2 = 0,82, \\ s_Y^2 &= \frac{313}{10} - 4,9^2 = 7,29, \\ s_{XY} &= \frac{345,2}{10} - 7,5 \cdot 4,9 = -2,23.\end{aligned}$$

Temos $b = -2,23/0,82 = -2,71$ e $a = 4,9 + 2,71 \cdot 7,5 = 25,25$ co que a ecuación da recta de regresión é

$$y = 25,25 - 2,71x.$$

A estimación do coeficiente de correlación (Subsección 6.1.2) é

$$r = \frac{-2,23}{\sqrt{0,82 \cdot 7,29}} = -0,91,$$

de xeito que a calidade da aproximación parece bastante boa.

A estimación do coeficiente de determinación (Sección 6.2) é $r^2 = 0,830$. Isto interprétase do seguinte xeito: o 83% da variabilidade da variable Y está explicada polo modelo de regresión.

Para contrasta-lo modelo de regresión linear temos que facer

$$H_0: \rho = 0, \quad H_1: \rho \neq 0.$$

Empregamos pois a técnica de análise da varianza, ANOVA (Subsección 6.2.1). Os datos necesarios están recollidos na seguinte táboa:

variabilidade	g.l.	SS	MS	cociente
regresión	1	$SS_R = 10 \cdot 0,83 \cdot 7,29 = 60,50$	$MS_R = 60,50$	39,02
erro	8	$SS_E = 10(1 - 0,83)7,29 = 12,40$	$MS_E = \frac{12,40}{8} = 1,55$	
total	9	$SS_Y = 10 \cdot 7,29 = 72,9$		

Como $P = P(F_{1,8} \geq 39,02) < 0,01$ é un número moi pequeno (de feito, empregando software estatístico temos $P = 0,00025$), rexeitámo-la hipótese nula. Concluimos que hai evidencia significativa de que o modelo de regresión linear é válido. \square

8.2. Exame de xuño de 2019

Problema 8.6. Moi recentemente, o xornal THE SUN publicou os resultados dun estudo sobre o peso dos paquetes de patacas fritas que as distintas cadeas de comida rápida serven en Inglaterra. O estudo consistiu en comprar tres paquetes de patacas de cada cadea en diferentes establecementos da mesma. En particular, para unha das cadeas, os resultados obtidos foron: 106g, 102g e 108g.

1. A partir da mostra, calcula un intervalo de confianza, cun nivel de confianza do 95 %, para o peso medio dos paquetes de patacas na devandita cadea.
2. Pódese afirmar, desde o punto de vista estatístico, que o peso medio real dos paquetes de patacas fritas nesa cadea é inferior a 108g?

Solución. Considerámo-la variable aleatoria X ="peso dun paquete de patacas fritas". Organizámo-los cálculos para obte-la media e cuasi-varianza mostral:

X	X^2
106	11236
102	10404
108	11664
Σ	316 33304

De aquí obtemos $n = 3$, $\bar{X} = \frac{316}{3} = 105,333$, $s_n^2 = \frac{33304}{3} - 105,333^2 = 6,222$, e así, $s_{n-1} = \sqrt{\frac{3}{2}} \cdot 6,222 = 3,055$.

Calculamos un intervalo de confianza para unha media (Subsección 2.2.3) empregando o estatístico

$$\frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}},$$

que segue unha distribución t_{n-1} . Despexando μ da desigualdade

$$\left| \frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}} \right| \leq t_{n-1, \alpha/2},$$

obtense a fórmula

$$\bar{X} \pm t_{n-1, \alpha/2} \frac{s_{n-1}}{\sqrt{n}}.$$

O nivel de significación é $\alpha = 0,05$. Calculamos $t_{2, 0,025} = 4,303$. Substituíndo na fórmula

$$105,333 \pm 4,303 \cdot \frac{3,055}{\sqrt{3}} = 105,333 \pm 7,589,$$

de onde se obtén o intervalo [97,744, 112,922].

Conclusión: cun nivel de confianza do 95.0 %, a media do peso dun paquete de patacas fritas atópase entre 97.744 e 112.922.

Agora facémo-lo contraste de hipóteses

$$H_0: \mu \geq 108, \quad H_1: \mu < 108.$$

Este é un contraste de hipóteses para unha media (Sección 3.1). Para iso empregámo-lo estatístico

$$\frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}},$$

que segue unha distribución t_{n-1} .

O valor no estatístico é

$$\frac{105,333 - 108}{3,055/\sqrt{3}} = -1,512.$$

Calculámo-lo valor P como $P = P(t_2 < -1,512) = 0,1349$, que é un valor relativamente grande.

Conclusión: Aceptamos H_0 , e concluímos que non hai evidencia significativa, ata un nivel de confianza do 86.5 %, de que a media de peso dun paquete de patacas fritas sexa menor ca 108. \square

Problema 8.7. Para saber se o olor a lavanda na sala de espera dos dentistas diminúe a ansiedade dos pacientes, un equipo de investigadores seleccionou a 597 pacientes que dividiu aleatoriamente en dous grupos. Os do primeiro grupo (310 pacientes), que chamaremos “grupo de control”, esperaron en salas sen aroma especial, mentres que os do segundo grupo (287 pacientes), que chamaremos “grupo de tratamento”, esperaron en salas con aroma a lavanda. Para determina-lo nivel de ansiedade, tódolos pacientes se someteron a diferentes test psicolóxicos que permiten medilo. Se nos test de ansiedade a media do grupo de control foi de 15.40 cunha cuasi-desviación típica de 4.18, e no grupo de tratamento a media mostral foi 11.74 cunha cuasi-desviación típica de 4.10, ¿podemos afirmar que o aroma de lavanda nas salas de espera dos dentistas axuda a reduci-lo nivel de ansiedade nos pacientes? NOTA: supoñede que as varianzas poboacionais son iguais.

Solución. Considerámo-las variables aleatorias X = “nivel de ansiedade no grupo de control” e Y = “nivel de ansiedade no grupo de tratamento”.

Temos $n_1 = 310$, $\bar{X} = 15,4$, $s_1 = 4,18$ e $n_2 = 287$, $\bar{Y} = 11,74$, $s_2 = 4,1$.

Asumimos que as varianzas das dúas poboacións son iguais.

Facémo-lo contraste de hipóteses

$$H_0: \mu_1 - \mu_2 \leq 0, \quad H_1: \mu_1 - \mu_2 > 0.$$

Este é un contraste de hipóteses para unha diferenza de medias (Subsección 4.1.2). Para iso empregámo-lo estatístico

$$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$$

que segue unha distribución $t_{n_1+n_2-2}$.

Aquí considerámo-la cuasi-varianza ponderada, que se define como

$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}.$$

Substituíndo na fórmula da cuasi-varianza ponderada obtemos

$$s_p = \sqrt{\frac{(310 - 1) \cdot 4,18^2 + (287 - 1) \cdot 4,1^2}{310 + 287 - 2}} = 4,142.$$

O valor no estatístico é

$$\frac{(15,4 - 11,74) - 0}{4,142 \sqrt{\frac{1}{310} + \frac{1}{287}}} = 10,788.$$

Calculámo-lo valor P como $P = P(t_{595} > 10,788) = 0,3 \cdot 10^{-24}$, que é un valor pequeno.

Conclusión: Rexeitamos H_0 , e concluimos que hai evidencia significativa, cun nivel de confianza do 99.9%, de que, en media, o nivel de ansiedade no grupo de control é maior có nivel de ansiedade no grupo de tratamento.

Por tanto, o aroma a lavanda na sala de espera dos dentistas axuda a reduci-lo nivel de ansiedade nos pacientes. \square

Problema 8.8. Para analiza-lo risco de sufrir un aborto espontáneo nos embarazos de mulleres hipertensas tratadas con inhibidores da encima convertidora de anxiotensina (IECA) durante o primeiro trimestre do embarazo, estudáronse 329 casos nos que se observaron 47 abortos espontáneos.

1. Se a taxa de abortos espontáneos na poboación fose do 10%, poderíase afirmar que o tratamento con IECA no primeiro trimestre de embarazo incrementa a porcentaxe de abortos espontáneos?
2. Cal tería que se-lo tamaño mostral mínimo para poder estimar, a un nivel de confianza do 95.5%, a proporción de abortos espontáneos na poboación cun erro inferior ó 2%?

Solución. Considerámo-la variable aleatoria X ="abortos espontáneos de mulleres hipertensas tratadas con IECA durante o primeiro trimestre do embarazo".

Temos $n = 329$, e $\hat{p} = 0,143$.

Facémo-lo contraste de hipóteses

$$H_0: p \leq 0,1, \quad H_1: p > 0,1.$$

Este é un contraste de hipóteses para unha proporción (Sección 3.3). Para iso empregámo-lo estatístico

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}},$$

que segue unha distribución normal estándar.

O valor no estatístico é

$$\frac{0,143 - 0,1}{\sqrt{\frac{0,1(1-0,1)}{329}}} = 2,591.$$

Calculámo-lo valor P como $P = P(Z > 2,591) = 0,0048$, que é un valor pequeno.

Conclusión: Rexeitamos H_0 , e concluimos que hai evidencia significativa, cun nivel de confianza do 99.5%, de que a proporción de abortos espontáneos de mulleres hipertensas tratadas con IECA durante o primeiro trimestre do embarazo é maior ca 10.0%.

Para estima-lo tamaño da mostra para unha proporción (Observación 2.25), empregámo-lo estatístico

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}},$$

que ten distribución normal estándar. Despexando p da desigualdade

$$\left| \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \right| \leq Z_{\alpha/2},$$

obtémo-la fórmula

$$\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}.$$

A estimación do erro é $Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$.

Neste case non temos unha estimación da proporción \hat{p} . É sinxelo ver que a función $x \mapsto \sqrt{x(1-x)}$ alcanza o seu máximo no intervalo $[0, 1]$ no punto $x = 1/2$. Por tanto, necesitamos despegar n da desigualdade $Z_{\alpha/2} \sqrt{\frac{0,5(1-0,5)}{n}} < \epsilon$, onde ϵ é o valor fixado polo problema. Así, obtense $n > \left(\frac{Z_{\alpha/2}}{2\epsilon}\right)^2$.

O nivel de significación é $\alpha = 0,045$. Mirando as táboas obtemos $Z_{0,0225} = 2,005$. Neste caso $\epsilon = 0,02$. Substituíndo na fórmula $n > \left(\frac{2,005}{2 \cdot 0,02}\right)^2 = 2511,65$.

Conclusión: para que a diferenza entre a proporción mostral e a proporción poboacional de abortos espontáneos de mulleres hipertensas tratadas con IECA durante o primeiro trimestre do embarazo sexa como moito de $\pm 0,02$ cun nivel de confianza do 95,5 %, teriamos que tomar unha mostra de polo menos 2512 elementos. \square

Problema 8.9. Co obxectivo de estudar a relación entre a aparición de depresión post-parto e o nivel de seguridade alimentaria, observáronse 325 casos de mulleres seleccionadas aleatoriamente en centros de saúde no oeste da cidade de Teherán (Irán). Clasifícase, de acordo coa seguridade alimentaria, ós fogares das devanditas mulleres en tres niveles: A1: Alimentación asegurada, A2: Alimentación non asegurada pero sen fame, A3: Alimentación non asegurada e con fame moderada ou severa. Dos 325 casos, 214 eran de fogares do tipo A1, 56 do tipo A2, e 55 do tipo A3. Dos 115 casos de depresión post-parto, 51 eran en mulleres con fogares de nivel A1, e 24 en mulleres con fogares de nivel A2.

1. Constrúe a táboa de continxencia e realiza o test estatístico adecuado para comprobar se hai relación entre a seguridade alimentaria no fogar e o feito de sufrir de depresión post-parto entre as mulleres da cidade de Teherán.
2. O test anterior, ¿é unha proba de independencia ou é unha proba de homoxeneidade? Razona a resposta.

Solución. Temos tres poboacións dependendo da seguridade alimentaria e a variable aleatoria Y ="depresión postparto". En primeiro lugar construímola táboa de continxencia:

Alimentación \ depresión	si	non	Σ
A1	51	163	214
A2	24	32	56
A3	40	15	55
Σ	115	210	325

Como o tamaño da mostra está determinado en toda a poboación, e o investigador simplemente clasifica os datos en dúas categorías, trátase dun contraste de independencia para datos categóricos (Sección 5.1). Por tanto, temos que face-lo contraste de hipóteses:

$$H_0: p_{ij} = p_i \cdot p_j, \quad i \in \{1, 2, 3\}, \quad j \in \{1, 2\}.$$

A continuación calculámo-las frecuencias esperadas, no suposto de que a hipótese nula sexa certa, mediante a fórmula $\widehat{E}_{ij} = \frac{n_{i.}n_{.j}}{n}$ (en verde), e tamén os valores intermedios do estatístico $(n_{ij} - \widehat{E}_{ij})^2 / \widehat{E}_{ij}$ (en vermello), obtendo:

Alimentación \ depresión	si	non	Σ
A1	51 75.72 8.07	163 138.28 4.42	214
A2	24 19.82 0.88	32 36.18 0.48	56
A3	40 19.46 21.67	15 35.54 11.87	55
Σ	115	210	325

Calcúlase o valor no estatístico, que consiste en suma-los valores vermellos. O resultado é 47.4.

O estatístico $\sum_{i,j} \frac{(n_{ij} - \widehat{E}_{ij})^2}{\widehat{E}_{ij}}$ segue unha distribución χ^2 con $(3 - 1)(2 - 1) = 2$ graos de liberdade. Calculando o valor P temos $P = P(\chi_2^2 \geq 47,4) = 0,5 \cdot 10^{-10}$.

Conclusión: rexeitámo-la hipótese nula, e por tanto, temos evidencia significativa, de que hai relación entre as dúas variables. \square

Problema 8.10. Co obxectivo de facer un modelo linear para predici-la altura dunha persoa a partir da lonxitude da súa tibia, nunha mostra aleatoria de 20 persoas médronse en centímetros tanto a súa tibia dereita (variable X), como a súa altura (variable Y) obténdose os seguintes valores:

$$\begin{aligned} \sum X &= 72,27; & \sum Y &= 322,48; \\ \sum X^2 &= 262,29; & \sum XY &= 1168,05; & \sum Y^2 &= 5206,53. \end{aligned}$$

1. Calcula a recta de regresión.
2. Calcula o coeficiente de determinación r^2 e interpreta o seu resultado.
3. Contrasta o modelo de regresión.

Nota: tomar 4 díxitos de precisión nos cálculos.

Solución. Considerámo-las seguintes variables aleatorias X ="lonxitude da tibia dereita" e Y ="altura".

Organizámo-los cálculos nunha táboa.

	X	Y	X^2	XY	Y^2
Σ	72.27	322.48	262.29	1168.05	5206.53

Temos $n = 20$ datos e

$$\begin{aligned}\bar{X} &= \frac{72,27}{20} = 3,613, \\ \bar{Y} &= \frac{322,48}{20} = 16,124, \\ s_X^2 &= \frac{262,29}{20} - 3,613^2 = 0,057, \\ s_Y^2 &= \frac{5206,53}{20} - 16,124^2 = 0,343, \\ s_{XY} &= \frac{1168,05}{20} - 3,613 \cdot 16,124 = 0,138.\end{aligned}$$

De aquí obtemos $b = 0,138/0,057 = 2,424$ e $a = 16,124 - 2,424 \cdot 3,613 = 7,367$, co que a ecuación da recta de regresión é

$$y = 7,367 + 2,424x.$$

A estimación do coeficiente de correlación é

$$r = \frac{0,138}{\sqrt{0,057 \cdot 0,343}} = 0,989.$$

A calidade da aproximación é forte.

O coeficiente de determinación vén dado por $r^2 = 0,978$. Isto interprétase do seguinte xeito: o 97,8% da variabilidade da variable Y está explicada polo modelo de regresión.

Contrastámo-la validez do modelo de regresión linear. Para iso facémo-lo contraste de hipóteses

$$H_0: \rho = 0, \quad H_1: \rho \neq 0.$$

Empregamos pois a técnica de análise da varianza, ANOVA. Os datos necesarios están recollidos na seguinte táboa:

variabilidade	g.l.	SS	MS	cociente
regresión	1	$SS_R = 20 \cdot 0,978 \cdot 0,343 = 6,71$	$MS_R = 6,71$	789,789
erro	18	$SS_E = 20 \cdot (1 - 0,978) \cdot 0,343 = 0,153$	$MS_E = \frac{0,153}{18} = 0,008$	
total	19	$SS_Y = 20 \cdot 0,343 = 6,862$		

Como $P = P(F_{1,18} \geq 789,789) = 0,3 \cdot 10^{-15}$ é un valor pequeno, rexeitámo-la hipótese nula. Concluimos que hai evidencia significativa de que o modelo de regresión linear é válido.

□

8.3. Exame de maio de 2021

Problema 8.11. Crese que as mellores condicións de vida e a desaparición de moitas enfermidades infecciosas levaron a unha aceleración do crecemento das poboacións dos países desenvolvidos. Para contrastalo, valorouse a altura X en cms de 127 adultos (homes) da poboación española en 2004 e comparáronse os resultados con estudos realizados antes de 1990, nos que a media poboacional era de 174,6cms. Obtívose unha media $\bar{X} = 177,33$ e unha cuasi-desviación típica $s = 3,26$.

1. ¿Apoia este estudo a idea de que a altura da poboación aumentou significativamente entre 1990 e 2004? Razoa o resultado usando o valor p .
2. Dá un intervalo de estimación para a media en 2004, cun nivel de confianza do 95%.
3. Se imos facer outro estudo e supoñemos que $\sigma = 3,26$ cms, ¿que tamaño de mostra necesitamos para que o erro de estimación da media sexa inferior a 1cm? Usar $\alpha = 5\%$.

Solución. Considerámo-la variable aleatoria X ="altura de adultos españois".

Temos $n = 127$, $\bar{X} = 177,33$, e $s_{n-1} = 3,26$.

Facémo-lo contraste de hipóteses

$$H_0: \mu \leq 174,6, \quad H_1: \mu > 174,6.$$

Este é un contraste de hipóteses para unha media (Sección 3.1). Para iso empregámo-lo estatístico

$$\frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}},$$

que segue unha distribución t -Student con $n - 1$ graos de liberdade.

O valor no estatístico é

$$\frac{177,33 - 174,6}{3,26/\sqrt{127}} = 9,437.$$

Mirámo-lo valor P na táboa para obter $P = P(t_{126} > 9,437) < 0,0005$ (De feito, $P = P(t_{126} > 9,437) = 0,1 \cdot 10^{-15}$), que é un valor moi pequeno.

Conclusión: Rexeitamos H_0 , e concluimos que hai evidencia significativa, cun nivel de confianza do 99,95%, de que a media de altura de adultos españois é maior ca 174,6.

Agora calculamos un intervalo de confianza para unha media (Subsección 2.2.3) empregando o estatístico

$$\frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}},$$

que segue unha distribución t -Student con $n - 1$ graos de liberdade. Despexando μ da desigualdade

$$\left| \frac{\bar{X} - \mu}{s_{n-1}/\sqrt{n}} \right| \leq t_{n-1, \alpha/2},$$

obtense a fórmula

$$\bar{X} \pm t_{n-1, \alpha/2} \frac{s_{n-1}}{\sqrt{n}}.$$

O nivel de significación é $\alpha = 0,05$. O valor en táboa máis próximo é $t_{100, 0,025} = 1,984$. Nesta solución, para da-lo resultado máis correcto posible, empregarémolo valor máis exacto $t_{126, 0,025} = 1,979$, pero o resultado é practicamente o mesmo. Substituímos na fórmula

$$177,33 \pm 1,979 \cdot \frac{3,26}{\sqrt{127}} = 177,33 \pm 0,572.$$

Por tanto obtense o intervalo $[176,758, 177,902]$.

Conclusión: cun nivel de confianza do 95,0%, a media de altura de adultos españois atópase entre 176,758 e 177,902.

Para estima-lo tamaño da mostra para unha media (Observación 2.5), empregámo-lo estatístico

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}},$$

que segue unha distribución normal estándar. Neste caso a varianza poboacional é coñecida. Despexando μ da desigualdade

$$\left| \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \right| \leq Z_{\alpha/2},$$

obtense a fórmula

$$\bar{X} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

A estimación do erro é

$$Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}.$$

Queremos $Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \epsilon$, onde ϵ é o valor fixado polo problema. Despexando obtense $n \geq \left(Z_{\alpha/2} \sigma / \epsilon \right)^2$.

O nivel de significación é $\alpha = 0,05$. Calculamos $Z_{0,025} = 1,96$. Neste caso $\epsilon = 1$. Substituíndo na fórmula, $n \geq (1,96 \cdot 3,26/1)^2 = 40,825$.

Conclusión: para que a diferenza entre a media mostral e a media poboacional de altura de adultos españois sexa como moito de ± 1 cun nivel de confianza do 95,0%, teriamos que tomar unha mostra de polo menos 41 elementos. \square

Problema 8.12. No mesmo estudo anterior, atopouse en 129 mulleres adultas que a altura media era de 163,96cms, con $s = 3,96$ cms. ¿Pódese afirmar que a altura media dos homes é maior cá das mulleres, cun nivel de significación de 10%?

Solución. Considerámo-las variables aleatorias X ="altura de homes españois" e Y ="altura de mulleres españolas".

Temos $n_1 = 127$, $\bar{X} = 177,33$, $s_1 = 3,26$, $n_2 = 129$, $\bar{Y} = 163,96$, $s_2 = 3,96$.

Asumimos que as varianzas das dúas poboacións son iguais.

Facémo-lo contraste de hipóteses

$$H_0: \mu_1 \leq \mu_2, \quad H_1: \mu_1 > \mu_2.$$

Este é un contraste de hipóteses para unha diferenza de medias (Subsección 4.1.2). Para iso empregámo-lo estatístico

$$\frac{(\bar{X} - \bar{Y}) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$$

que segue unha distribución t -Student con $n_1 + n_2 - 2$ graos de liberdade.

Aquí a cuasi-varianza mostral conxunta é

$$s_p^2 = \frac{(n_1 - 1) s_1^2 + (n_2 - 1) s_2^2}{n_1 + n_2 - 2}.$$

Substituíndo na fórmula da cuasi-varianza mostral conxunta:

$$s_p^2 = \frac{(127 - 1) 3,26^2 + (129 - 1) 3,96^2}{127 + 129 - 2} = 13,174,$$

polo que $s_p = 3,63$.

O valor no estatístico é

$$\frac{(177,33 - 163,96) - 0}{3,63 \sqrt{\frac{1}{127} + \frac{1}{129}}} = 29,467.$$

O nivel de significación é $\alpha = 0,1$. O valor en táboa máis próximo é $t_{200,0,1} = 1,286$. Con software informático obtense un valor un pouco máis exacto: $t_{254,0,1} = 1,285$. En calquera caso, $29,467 \notin (-\infty, 1,285]$.

Conclusión: Rexeitamos H_0 , e concluimos que hai evidencia significativa, cun nivel de confianza do 90,0%, de que a altura media de homes españois é maior cá das mulleres. \square

Problema 8.13. Sábese que o uso prolongado de antibióticos pode causar toxicidade neurolóxica. Cando 10000 persoas recibiron metronidazol, 11 delas sufriron ataxia (movemento muscular non coordinado).

1. ¿Pódese dicir que este medicamento causa ataxia en máis dun caso por cada 1000 pacientes?
2. Indica un intervalo de estimación da frecuencia de aparición de síntomas de ataxia.

Solución. Considerámo-la variable aleatoria X ="casos de ataxia en persoas que recibiron metronidazol".

Temos $n = 10000$, $\hat{p} = 0,0011$.

Facémo-lo contraste de hipóteses

$$H_0: p \leq 0,001, \quad H_1: p > 0,001.$$

Este é un contraste de hipóteses para unha proporción (Sección 3.3). Para iso empregámo-lo estatístico

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}},$$

que segue unha distribución normal estándar.

O valor no estatístico é

$$\frac{0,0011 - 0,001}{\sqrt{\frac{0,001(1-0,001)}{10000}}} = 0,3164.$$

Calculámo-lo valor P (aproximadamente), mirando na táboa, como $P = P(Z > 0,32) = 0,374$, que é un valor relativamente grande.

Conclusión: Aceptamos H_0 , e concluimos que non hai evidencia significativa de que a proporción de casos de ataxia en persoas que recibiron metronidazol sexa maior có 0,1%.

Calculamos agora un intervalo de confianza para unha proporción (Sección 2.4) empregando o estatístico

$$\frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}},$$

que segue unha distribución normal estándar. Despejando p da desigualdade

$$\left| \frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \right| \leq Z_{\alpha/2},$$

obtense a fórmula

$$\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}.$$

O nivel de significación é $\alpha = 0,1$. Calculamos $Z_{0,05} = 1,6449$. Substituímos na fórmula

$$0,0011 \pm 1,6449 \sqrt{\frac{0,0011(1-0,0011)}{10000}} = 0,0011 \pm 0,0005.$$

Por tanto obtense o intervalo $[0,0006, 0,0016]$.

Conclusión: cun nivel de confianza do 90,0%, a proporción de casos de ataxia en persoas que recibiron metronidazol atópase entre 0,06% e 0,16%. \square

Problema 8.14. En 1977 levouse a cabo un experimento para coñecer a incidencia de efectos secundarios asociados ó uso de minociclina, un antibiótico prescrito para tratar a acné. Foron incluídos 45 pacientes no grupo de tratamento, e 44 pacientes no grupo placebo. Dos pacientes do grupo de tratamento, 33 presentaron síntomas vestibulares (sensación de vertixe), en comparación con 4 persoas que tiveron síntomas no grupo placebo. Elabora a táboa de continxencia axeitada e determina se existe unha asociación entre ter seguido o tratamento e sufrir síntomas vestibulares.

Solución. Temos 2 poboacións, dependendo do “grupo” (paciente ou placebo), e a variable aleatoria Y = “síntomas”.

En primeiro lugar construímo-la táboa de continxencia:

grupo \ síntomas	si	non	tamaño
paciente	33	12	45
placebo	4	40	44
Σ	37	52	89

Temos que face-lo contraste de hipóteses:

$$H_0: p_{11} = p_{21}, p_{12} = p_{22}.$$

Este é un contraste de hipóteses para homoxeneidade de datos categóricos (Sección 5.2), xa que o tamaño da mostra en cada poboación é fixado polo investigador. Para iso empregámo-lo estatístico

$$\sum_{i,j} \frac{(n_{ij} - \widehat{E}_{ij})^2}{\widehat{E}_{ij}},$$

que segue unha distribución χ^2 de Pearson con $(f-1)(c-1)$ graos de liberdade.

O número de graos de liberdade da distribución é $(2-1)(2-1) = 1$.

A continuación calculámo-las frecuencias esperadas, no suposto de que a hipótese nula sexa certa, mediante a fórmula $\widehat{E}_{ij} = \frac{n_{i.}n_{.j}}{n}$:

grupo \ síntomas	si	non	tamaño
paciente	18.71	26.29	45
placebo	18.29	25.71	44
Σ	37	52	89

Agora calculámo-los valores intermedios do estatístico $(n_{ij} - \widehat{E}_{ij})^2 / \widehat{E}_{ij}$:

grupo \ síntomas	si	non	Σ
paciente	10.919	7.769	
placebo	11.167	7.946	
Σ			37.8

A suma dos valores intermedios, que coincide co valor no estatístico, é 37.8.

Calculámo-lo valor P mirando as táboas para obter $P = P(\chi_1^2 > 37,8) < 0,001$. En realidade, con software informático obtense $P = P(\chi_1^2 > 37,8) = 0,8 \cdot 10^{-9}$, que en todo caso é un valor moi pequeno.

Conclusión: Rexeitamos H_0 , e concluímos que hai evidencia significativa, cun nivel de confianza do 99,9%, de que hai relación entre ter seguido o tratamento e sufrir síntomas vestibulares. \square

Problema 8.15. A ataxia cerebelosa caracterízase por dificultades no equilibrio. Estase a desenvolver unha aplicación para teléfono móbil capaz de medir, mediante un acelerador situado á altura do esterno, o equilibrio estático e dinámico en pacientes con esa enfermidade. Para estima-la validez da aplicación, as medicións Y obtidas con ela relacionáronse co “índice de estabilidade postural estática” X obtido mediante unha plataforma colocada no chan. En 6 pacientes obtivéronse os seguintes resultados:

X	45.0	61.7	129.0	392.0	285.5	209.6
Y	3.76	4.04	4.34	4.48	4.80	4.12

Pídese:

1. Estima-la puntuación Y que obtería na nova aplicación un paciente con $X = 515,0$.
2. Calcular e interpreta-lo coeficiente de determinación r^2 .
3. Realiza-lo contraste ANOVA para a regresión lineal.

Solución. Considerámo-las variables aleatorias X = “índice de estabilidade postural estática” e Y = “equilibrio estático e dinámico en pacientes con ataxia”.

Organizámo-los cálculos nunha táboa.

X	Y	X^2	XY	Y^2
45.0	3.76	2025.0	169.2	14.138
61.7	4.04	3806.89	249.268	16.322
129.0	4.34	16641.0	559.86	18.836
392.0	4.48	153664.0	1756.16	20.07
285.5	4.8	81510.25	1370.4	23.04
209.6	4.12	43932.16	863.552	16.974
Σ	1122.8	301579.3	4968.44	109.38

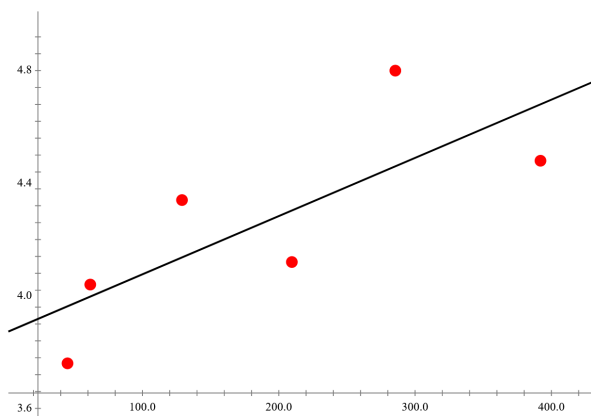


Figura 8.1: Os puntos e a súa recta de regresión

Temos $n = 6$ datos e

$$\begin{aligned}\bar{X} &= \frac{1122,8}{6} = 187,133, \\ \bar{Y} &= \frac{25,54}{6} = 4,257, \\ s_X^2 &= \frac{301579,3}{6} - 187,133^2 = 15244,332, \\ s_Y^2 &= \frac{109,38}{6} - 4,257^2 = 0,111, \\ s_{XY} &= \frac{4968,44}{6} - 187,133 \cdot 4,257 = 31,509.\end{aligned}$$

De aquí obtemos

$$\begin{aligned}b &= 31,509 / 15244,332 = 0,002, \\ a &= 4,257 - 0,002 \cdot 187,133 = 3,87,\end{aligned}$$

co que a ecuación da recta de regresión é

$$y = 3,87 + 0,002x.$$

Avaliando na recta de regresión, para “índice de estabilidade postural estática” $x = 515,0$ estímase o “equilibrio estático e dinámico en pacientes con ataxia”

$$y = 3,87 + 0,002 \cdot 515,0 = 4,934.$$

A estimación do coeficiente de correlación é

$$r = \frac{31,509}{\sqrt{15244,332 \cdot 0,111}} = 0,767.$$

A calidade da aproximación é moderada.

O coeficiente de determinación vén dado por $r^2 = 0,588$. Isto interprétase do seguinte xeito: o 58,8% da variabilidade da variable Y está explicada polo modelo de regresión.

Comprobámo-la validez do modelo de regresión linear.
 Facémo-lo contraste de hipóteses

$$H_0: \rho = 0, \quad H_1: \rho \neq 0.$$

Este é un contraste de hipóteses para a validez do modelo de regresión linear (Subsección 6.2.1). Para iso empregámo-lo estatístico que se obtén despois de dispoñer dos cálculos nunha táboa ANOVA e que segue unha distribución F de Snedecor con $(1, n - 2)$ graos de liberdade.

	g.l.	SS	MS	cociente
regresión	1	$SS_R = 6 \cdot 0,588 \cdot 0,111 = 0,391$	$MS_R = 0,391$	5,714
erro	4	$SS_E = 6 \cdot (1 - 0,588) \cdot 0,111 = 0,274$	$MS_E = \frac{0,274}{4} = 0,068$	
total	5	$SS_Y = 6 \cdot 0,111 = 0,664$		

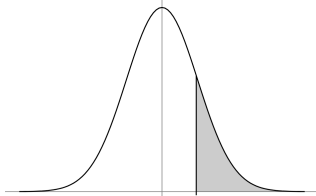
Temos que mirar en dúas táboas da F de Snedecor para estima-lo valor P e chegar a que, se $P = P(F_{1,4} > 5,714)$, entón $0,05 < P < 0,1$. Calculando o valor P con software estatístico obtense $P = P(F_{1,4} > 5,714) = 0,0751$, que é un valor relativamente pequeno.

Conclusión: Rexeitamos H_0 , e concluímos que hai evidencia significativa, cun nivel de confianza do 90%, de que o modelo de regresión linear é válido. Non obstante existen dúbidas, xa que cun nivel de confianza do 95% non teríamos evidencia estatística da súa validez. \square

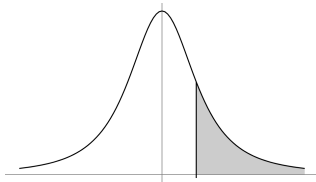
Capítulo 9

Táboas estatísticas

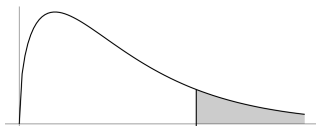
Táboas estatísticas que se empregan neste curso.



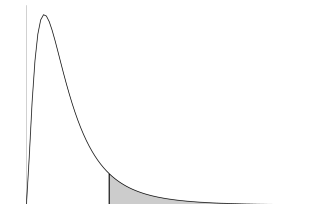
Táboa da distribución normal



Táboa da distribución t de Student



Táboas da distribución χ^2 de Pearson

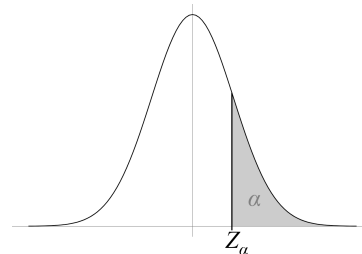


Táboas da distribución F de Fisher-Snedecor

Distribución normal tipificada

Área de cola derecha

$$\int_{Z_\alpha}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx = \alpha$$

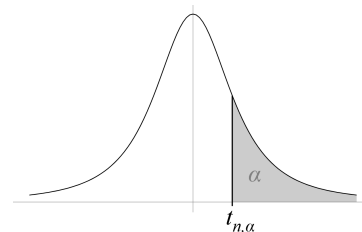


Z_α	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.50000	0.49601	0.49202	0.48803	0.48405	0.48006	0.47608	0.47210	0.46812	0.46414
0.1	0.46017	0.45620	0.45224	0.44828	0.44433	0.44038	0.43644	0.43251	0.42858	0.42465
0.2	0.42074	0.41683	0.41294	0.40905	0.40517	0.40129	0.39743	0.39358	0.38974	0.38591
0.3	0.38209	0.37828	0.37448	0.37070	0.36693	0.36317	0.35942	0.35569	0.35197	0.34827
0.4	0.34458	0.34090	0.33724	0.33360	0.32997	0.32636	0.32276	0.31918	0.31561	0.31207
0.5	0.30854	0.30503	0.30153	0.29806	0.29460	0.29116	0.28774	0.28434	0.28096	0.27760
0.6	0.27425	0.27093	0.26763	0.26435	0.26109	0.25785	0.25463	0.25143	0.24825	0.24510
0.7	0.24196	0.23885	0.23576	0.23270	0.22965	0.22663	0.22363	0.22065	0.21770	0.21476
0.8	0.21186	0.20897	0.20611	0.20327	0.20045	0.19766	0.19489	0.19215	0.18943	0.18673
0.9	0.18406	0.18141	0.17879	0.17619	0.17361	0.17106	0.16853	0.16602	0.16354	0.16109
1.0	0.15866	0.15625	0.15386	0.15151	0.14917	0.14686	0.14457	0.14231	0.14007	0.13786
1.1	0.13567	0.13350	0.13136	0.12924	0.12714	0.12507	0.12302	0.12100	0.11900	0.11702
1.2	0.11507	0.11314	0.11123	0.10935	0.10749	0.10565	0.10383	0.10204	0.10027	0.09853
1.3	0.09680	0.09510	0.09342	0.09176	0.09012	0.08851	0.08691	0.08534	0.08379	0.08226
1.4	0.08076	0.07927	0.07780	0.07636	0.07493	0.07353	0.07215	0.07078	0.06944	0.06811
1.5	0.06681	0.06552	0.06426	0.06301	0.06178	0.06057	0.05938	0.05821	0.05705	0.05592
1.6	0.05480	0.05370	0.05262	0.05155	0.05050	0.04947	0.04846	0.04746	0.04648	0.04551
1.7	0.04457	0.04363	0.04272	0.04182	0.04093	0.04006	0.03920	0.03836	0.03754	0.03673
1.8	0.03593	0.03515	0.03438	0.03362	0.03288	0.03216	0.03144	0.03074	0.03005	0.02938
1.9	0.02872	0.02807	0.02743	0.02680	0.02619	0.02559	0.02500	0.02442	0.02385	0.02330
2.0	0.02275	0.02222	0.02169	0.02118	0.02068	0.02018	0.01970	0.01923	0.01876	0.01831
2.1	0.01786	0.01743	0.01700	0.01659	0.01618	0.01578	0.01539	0.01500	0.01463	0.01426
2.2	0.01390	0.01355	0.01321	0.01287	0.01255	0.01222	0.01191	0.01160	0.01130	0.01101
2.3	0.01072	0.01044	0.01017	0.00990	0.00964	0.00939	0.00914	0.00889	0.00866	0.00842
2.4	0.00820	0.00798	0.00776	0.00755	0.00734	0.00714	0.00695	0.00676	0.00657	0.00639
2.5	0.00621	0.00604	0.00587	0.00570	0.00554	0.00539	0.00523	0.00508	0.00494	0.00480
2.6	0.00466	0.00453	0.00440	0.00427	0.00415	0.00402	0.00391	0.00379	0.00368	0.00357
2.7	0.00347	0.00336	0.00326	0.00317	0.00307	0.00298	0.00289	0.00280	0.00272	0.00264
2.8	0.00256	0.00248	0.00240	0.00233	0.00226	0.00219	0.00212	0.00205	0.00199	0.00193
2.9	0.00187	0.00181	0.00175	0.00169	0.00164	0.00159	0.00154	0.00149	0.00144	0.00139
3.0	0.00135	0.00131	0.00126	0.00122	0.00118	0.00114	0.00111	0.00107	0.00104	0.00100
Z_α	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09

Distribución t de Student

Abscisas $t_{n,\alpha}$ que deixan á súa dereita un área α
 nunha t -Student con n graos de liberdade

$$\int_{t_{n,\alpha}}^{\infty} c_n \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}} dx = \alpha$$

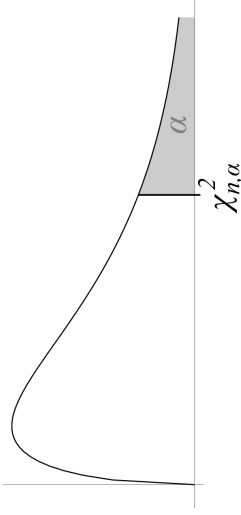


n \ \alpha	0.4000	0.3000	0.2000	0.1000	0.0500	0.0250	0.0100	0.0050	0.0010	0.0005
1	0.3249	0.7265	1.3764	3.0777	6.3138	12.706	31.820	63.656	318.30	636.61
2	0.2887	0.6172	1.0607	1.8856	2.9200	4.3027	6.9646	9.9248	22.327	31.599
3	0.2767	0.5844	0.9785	1.6377	2.3534	3.1824	4.5407	5.8409	10.214	12.924
4	0.2707	0.5686	0.9410	1.5332	2.1318	2.7764	3.7469	4.6041	7.1732	8.6103
5	0.2672	0.5594	0.9195	1.4759	2.0150	2.5706	3.3649	4.0321	5.8934	6.8688
6	0.2648	0.5534	0.9057	1.4398	1.9432	2.4469	3.1427	3.7074	5.2076	5.9588
7	0.2632	0.5491	0.8960	1.4149	1.8946	2.3646	2.9980	3.4995	4.7853	5.4079
8	0.2619	0.5459	0.8889	1.3968	1.8595	2.3060	2.8965	3.3554	4.5008	5.0413
9	0.2610	0.5435	0.8834	1.3830	1.8331	2.2622	2.8214	3.2498	4.2968	4.7809
10	0.2602	0.5415	0.8791	1.3722	1.8125	2.2281	2.7638	3.1693	4.1437	4.5869
11	0.2596	0.5399	0.8755	1.3634	1.7959	2.2010	2.7181	3.1058	4.0247	4.4370
12	0.2590	0.5386	0.8726	1.3562	1.7823	2.1788	2.6810	3.0545	3.9296	4.3178
13	0.2586	0.5375	0.8702	1.3502	1.7709	2.1604	2.6503	3.0123	3.8520	4.2208
14	0.2582	0.5366	0.8681	1.3450	1.7613	2.1448	2.6245	2.9768	3.7874	4.1405
15	0.2579	0.5357	0.8662	1.3406	1.7531	2.1314	2.6025	2.9467	3.7328	4.0728
16	0.2576	0.5350	0.8647	1.3368	1.7459	2.1199	2.5835	2.9208	3.6862	4.0150
17	0.2573	0.5344	0.8633	1.3334	1.7396	2.1098	2.5669	2.8982	3.6458	3.9651
18	0.2571	0.5338	0.8620	1.3304	1.7341	2.1009	2.5524	2.8784	3.6105	3.9216
19	0.2569	0.5333	0.8610	1.3277	1.7291	2.0930	2.5395	2.8609	3.5794	3.8834
20	0.2567	0.5329	0.8600	1.3253	1.7247	2.0860	2.5280	2.8453	3.5518	3.8495
21	0.2566	0.5325	0.8591	1.3232	1.7207	2.0796	2.5176	2.8314	3.5272	3.8193
22	0.2564	0.5321	0.8583	1.3212	1.7171	2.0739	2.5083	2.8188	3.5050	3.7921
23	0.2563	0.5317	0.8575	1.3195	1.7139	2.0687	2.4999	2.8073	3.4850	3.7676
24	0.2562	0.5314	0.8569	1.3178	1.7109	2.0639	2.4922	2.7969	3.4668	3.7454
25	0.2561	0.5312	0.8562	1.3163	1.7081	2.0595	2.4851	2.7874	3.4502	3.7251
26	0.2560	0.5309	0.8557	1.3150	1.7056	2.0555	2.4786	2.7787	3.4350	3.7066
27	0.2559	0.5306	0.8551	1.3137	1.7033	2.0518	2.4727	2.7707	3.4210	3.6896
28	0.2558	0.5304	0.8546	1.3125	1.7011	2.0484	2.4671	2.7633	3.4082	3.6739
29	0.2557	0.5302	0.8542	1.3114	1.6991	2.0452	2.4620	2.7564	3.3962	3.6594
30	0.2556	0.5300	0.8538	1.3104	1.6973	2.0423	2.4573	2.7500	3.3852	3.6460
40	0.2550	0.5286	0.8507	1.3031	1.6839	2.0211	2.4233	2.7045	3.3069	3.5510
50	0.2547	0.5278	0.8489	1.2987	1.6759	2.0086	2.4033	2.6778	3.2614	3.4960
60	0.2545	0.5272	0.8477	1.2958	1.6706	2.0003	2.3901	2.6603	3.2317	3.4602
80	0.2542	0.5265	0.8461	1.2922	1.6641	1.9901	2.3739	2.6387	3.1953	3.4163
100	0.2540	0.5261	0.8452	1.2901	1.6602	1.9840	2.3642	2.6259	3.1737	3.3905
200	0.2537	0.5252	0.8434	1.2858	1.6525	1.9719	2.3451	2.6006	3.1315	3.3398
500	0.2535	0.5247	0.8423	1.2832	1.6479	1.9647	2.3338	2.5857	3.1066	3.3101
∞	0.2533	0.5244	0.8416	1.2816	1.6449	1.9600	2.3263	2.5758	3.0902	3.2905
n / \alpha	0.4000	0.3000	0.2000	0.1000	0.0500	0.0250	0.0100	0.0050	0.0010	0.0005

Distribución χ^2 de Pearson

Abscisas $\chi_{n,\alpha}^2$ que deixan á súa dereita un área α nunha χ^2 con n graos de liberdade

$$\int_{\chi_{n,\alpha}^2}^{\infty} c_n x^{n/2-1} e^{-x/2} dx = \alpha$$

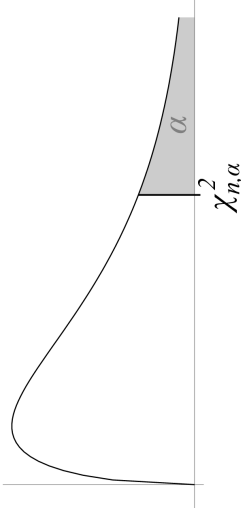


n\alpha	0.001	0.005	0.010	0.025	0.030	0.040	0.050	0.100	0.150	0.200	0.250	0.300	0.350	0.400
1	10.827	7.8794	6.6349	5.0239	4.7093	4.2179	3.8415	2.7055	2.0723	1.6424	1.3233	1.0742	0.8735	0.7083
2	13.815	10.596	9.2103	7.3778	7.0131	6.4378	5.9915	4.6052	3.7942	3.2189	2.7726	2.4079	2.0996	1.8326
3	16.266	12.838	11.344	9.3484	8.9473	8.3112	7.8147	6.2514	5.3170	4.6416	4.1083	3.6649	3.2831	2.9462
4	18.466	14.860	13.276	11.143	10.711	10.025	9.4877	7.7794	6.7449	5.9886	5.3853	4.8784	4.4377	4.0446
5	20.515	16.749	15.086	12.832	12.374	11.644	11.070	9.2364	8.1152	7.2893	6.6257	6.0644	5.5731	5.1319
6	22.457	18.547	16.811	14.449	13.967	13.197	12.591	10.644	9.4461	8.5581	7.8408	7.2311	6.6948	6.2108
7	24.321	20.277	18.475	16.012	15.509	14.703	14.067	12.017	10.747	9.8032	9.0371	8.3834	7.8061	7.2832
8	26.124	21.955	20.090	17.534	17.010	16.170	15.507	13.361	12.027	11.030	10.218	9.5245	8.9094	8.3505
9	27.877	23.589	21.666	19.022	18.479	17.608	16.919	14.683	13.288	12.242	11.388	10.656	10.006	9.4136
10	29.588	25.188	23.209	20.483	19.921	19.020	18.307	15.987	14.533	13.442	12.548	11.780	11.097	10.473
11	31.264	26.756	24.725	21.920	21.341	20.412	19.675	17.275	15.767	14.631	13.700	12.898	12.183	11.529
12	32.909	28.299	26.217	23.336	22.741	21.785	21.026	18.549	16.989	15.812	14.845	14.011	13.266	12.583
13	34.528	29.819	27.688	24.735	24.124	23.142	22.362	19.811	18.202	16.984	15.983	15.118	14.345	13.635
14	36.123	31.319	29.141	26.118	25.493	24.485	23.684	19.406	18.150	17.116	16.222	15.420	14.685	14.685
15	37.697	32.801	30.577	27.488	26.847	25.816	24.995	20.603	19.310	18.245	17.321	16.494	15.733	15.733
16	39.252	34.267	31.999	28.845	28.190	27.135	26.296	23.541	21.793	20.465	19.368	18.417	17.564	16.779
17	40.790	35.718	33.408	30.191	29.522	28.445	27.587	24.769	22.977	21.614	20.488	19.511	18.633	17.824
18	42.312	37.156	34.805	31.526	30.844	29.745	28.869	25.989	24.155	22.759	21.604	20.601	19.699	18.867
19	43.820	38.582	36.190	32.852	32.157	31.036	30.143	27.203	25.328	23.900	22.717	21.689	20.763	19.910
20	45.314	39.996	37.566	34.169	33.462	32.320	31.410	28.412	26.497	25.037	23.827	22.774	21.826	20.951
21	46.797	41.401	38.932	35.478	34.759	33.597	32.670	29.615	27.662	26.171	24.934	23.857	22.887	21.991
22	48.267	42.795	40.289	36.780	36.049	34.867	33.924	30.813	28.822	27.301	26.039	24.939	23.947	23.030
23	49.728	44.181	41.638	38.075	37.332	36.131	35.172	32.006	29.979	28.428	27.141	26.018	25.005	24.068
24	51.178	45.558	42.979	39.364	38.609	37.389	36.415	33.196	31.132	29.553	28.241	27.096	26.062	25.106
25	52.619	46.927	44.314	40.646	39.880	38.641	37.652	34.381	32.282	30.675	29.338	28.171	27.118	26.143
26	54.052	48.289	45.641	41.923	41.146	39.889	38.885	35.563	33.429	31.794	30.434	29.246	28.173	27.178
27	55.476	49.644	46.962	43.194	42.406	41.131	40.113	36.741	34.573	32.911	31.528	30.319	29.226	28.214
28	56.892	50.993	48.278	44.460	43.662	42.369	41.337	37.915	35.715	34.026	32.620	31.390	30.279	29.248
29	58.301	52.335	49.587	45.722	44.913	43.603	42.557	39.087	36.853	35.139	33.710	32.461	31.330	30.282
30	59.703	53.672	50.892	46.979	46.159	44.833	43.773	40.256	37.990	36.250	34.799	33.530	32.381	31.315
31	61.098	55.002	52.191	48.231	47.402	46.059	44.985	41.421	39.124	37.359	35.887	34.598	33.431	32.348
32	62.487	56.328	53.485	49.480	48.641	47.281	46.194	42.584	40.256	38.466	36.973	35.664	34.480	33.380
33	63.870	57.648	54.775	50.725	49.875	48.500	47.399	43.745	41.386	39.571	38.057	36.730	35.528	34.412
34	65.247	58.963	56.060	51.966	51.107	49.715	48.602	44.903	42.514	40.675	39.140	37.795	36.576	35.443
35	66.618	60.274	57.342	53.203	52.335	50.928	49.801	46.058	43.639	41.778	40.222	38.859	37.623	36.474
40	73.402	66.766	63.690	59.341	58.427	56.945	55.758	51.805	49.243	47.268	45.616	44.164	42.847	41.622
60	99.607	91.951	88.379	83.297	82.225	80.482	79.081	74.397	71.341	68.972	66.981	65.226	63.627	62.134
80	124.83	116.32	112.32	106.62	105.42	103.45	101.87	96.578	93.105	90.405	88.130	86.119	84.284	82.566
90	137.20	128.29	124.11	118.13	116.86	114.80	113.14	107.56	103.90	101.05	98.649	96.523	94.580	92.761
100	149.44	140.16	135.80	129.56	128.23	126.07	124.34	118.49	114.65	111.66	109.14	106.90	104.86	102.94
120	173.61	163.64	158.95	152.21	150.78	148.44	146.56	140.23	136.06	132.80	130.05	127.61	125.38	123.28
140	197.45	186.84	181.84	174.64	173.11	170.62	168.61	161.82	157.35	153.85	150.89	148.26	145.86	143.60
n/α	0.001	0.005	0.010	0.025	0.030	0.040	0.050	0.100	0.150	0.200	0.250	0.300	0.350	0.400

Distribución χ^2 de Pearson

Abscisas $\chi_{n,\alpha}^2$ que deixan á súa dereita un área α nunha χ^2 con n graos de liberdade

$$\int_{\chi_{n,\alpha}^2}^{\infty} c_n x^{n/2-1} e^{-x/2} dx = \alpha$$



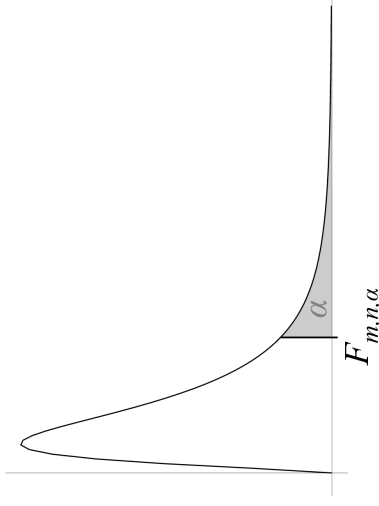
n\α	0.450	0.500	0.550	0.600	0.650	0.700	0.750	0.800	0.850	0.900	0.950	0.975	0.980	0.990	0.995
1	0.5707	0.4549	0.3573	0.2750	0.2059	0.1485	0.1015	0.0642	0.0358	0.0158	0.0039	0.0010	0.0006	0.0002	0.0000
2	1.5970	1.3863	1.1957	1.0217	0.8616	0.7133	0.5754	0.4463	0.3250	0.2107	0.1026	0.0506	0.0404	0.0201	0.0100
3	2.6430	2.3660	2.1095	1.8692	1.6416	1.4237	1.2125	1.0052	0.7978	0.5844	0.3518	0.2158	0.1848	0.1148	0.0717
4	3.6871	3.3567	3.0469	2.7528	2.4701	2.1947	1.9226	1.6488	1.3665	1.0636	0.7107	0.4844	0.4294	0.2971	0.2070
5	4.7278	4.3515	3.9959	3.6555	3.3251	2.9999	2.6746	2.3425	1.9938	1.6103	1.1455	0.8312	0.7519	0.5543	0.4117
6	5.7652	5.3481	4.9519	4.5702	4.1973	3.8276	3.4546	3.0701	2.6613	2.2041	1.6354	1.2373	1.1344	0.8721	0.6757
7	6.8000	6.3458	5.9125	5.4932	5.0816	4.6713	4.2549	3.8223	3.3583	2.8331	2.1673	1.6899	1.5643	1.2390	0.9893
8	7.8325	7.3441	6.8766	6.4226	5.9753	5.5274	5.0706	4.5936	4.0782	3.4895	2.7326	2.1797	2.0325	1.6465	1.3444
9	8.8632	8.3428	7.8434	7.3570	6.8763	6.3933	5.8988	5.3801	4.8165	4.1682	3.3251	2.7004	2.5324	2.0879	1.7349
10	9.8922	9.3418	8.8124	8.2955	7.7832	7.2672	6.7372	6.1791	5.5701	4.8652	3.9403	3.2470	3.0591	2.5582	2.1559
11	10.919	10.341	9.7831	9.2373	8.6952	8.1479	7.5841	6.9887	6.3364	5.5778	4.5748	3.8157	3.6087	3.0535	2.6032
12	11.946	11.340	10.755	10.182	9.6115	9.0343	8.4384	7.8073	7.1138	6.3038	5.2260	4.4038	4.1783	3.5706	3.0738
13	12.971	12.339	11.728	11.129	10.531	9.9257	9.2991	8.6339	7.9008	7.0415	5.8919	5.0088	4.7654	4.1069	3.5650
14	13.996	13.339	12.703	12.078	11.454	10.821	10.165	9.4673	8.6963	7.7895	6.5706	5.6287	5.3682	4.6604	4.0747
15	15.019	14.338	13.679	13.029	12.380	11.721	11.036	10.307	9.4993	8.5468	7.2609	6.2621	5.9849	5.2293	4.6009
16	16.042	15.338	14.655	13.982	13.309	12.624	11.912	11.152	10.309	9.3122	7.9616	6.9077	6.6142	5.8122	5.1422
17	17.064	16.338	15.632	14.937	14.240	13.530	12.791	12.002	11.124	10.085	8.6718	7.5642	7.2550	6.4078	5.6972
18	18.086	17.337	16.610	15.893	15.173	14.439	13.675	12.857	11.946	10.864	9.3905	8.2907	7.9062	7.0149	6.2648
19	19.106	18.337	17.589	16.850	16.108	15.351	14.562	13.715	12.772	11.650	10.117	8.9065	8.5670	7.6327	6.8440
20	20.127	19.337	18.568	17.808	17.045	16.265	15.451	14.578	13.603	12.442	10.850	9.5908	9.2367	8.2604	7.4338
21	21.147	20.337	19.548	18.768	17.984	17.182	16.344	15.444	14.439	13.239	11.591	10.282	9.9146	8.8972	8.0337
22	22.166	21.337	20.528	19.728	18.924	18.100	17.239	16.314	15.278	14.041	12.338	10.982	10.600	9.5425	8.6427
23	23.185	22.336	21.509	20.690	19.865	19.021	18.137	17.186	16.121	14.848	13.090	11.688	11.292	10.195	9.2604
24	24.203	23.336	22.490	21.652	20.808	19.943	19.037	18.061	16.968	15.658	13.848	12.401	11.991	10.856	9.8862
25	25.221	24.336	23.472	22.615	21.752	20.867	19.939	18.939	17.818	16.473	14.611	13.119	12.697	11.524	10.519
26	26.239	25.336	24.454	23.579	22.697	21.792	20.843	19.820	18.671	17.291	15.379	13.843	13.408	12.198	11.160
27	27.256	26.336	25.436	24.544	23.643	22.719	21.749	20.703	19.527	18.113	16.151	14.573	14.125	12.878	11.807
28	28.274	27.336	26.419	25.509	24.590	23.647	22.657	21.588	20.385	18.939	16.927	15.307	14.847	13.564	12.461
29	29.290	28.336	27.402	26.475	25.539	24.577	23.566	22.475	21.246	19.767	17.708	16.047	15.574	14.256	13.121
30	30.307	29.336	28.385	27.441	26.488	25.507	24.477	23.364	22.110	20.599	18.492	16.790	16.306	14.953	13.786
31	31.323	30.335	29.369	28.408	27.438	26.439	25.390	24.255	22.976	21.433	19.280	17.538	17.042	15.655	14.457
32	32.339	31.335	30.353	29.376	28.388	27.372	26.304	25.147	23.844	22.270	20.071	18.290	17.782	16.362	15.134
33	33.355	32.335	31.337	30.344	29.340	28.306	27.219	26.042	24.714	23.110	20.866	19.046	18.527	17.073	15.815
34	34.370	33.335	32.321	31.313	30.292	29.242	28.136	26.938	25.586	23.952	21.664	19.806	19.275	17.789	16.501
35	35.385	34.335	33.306	32.282	31.245	30.178	29.054	27.835	26.460	24.796	22.465	20.569	20.027	18.508	17.191
40	40.458	39.335	38.232	37.134	36.020	34.871	33.660	32.345	30.856	29.050	26.509	24.433	23.837	22.164	20.706
60	60.712	59.334	57.977	56.620	55.239	53.809	52.293	50.640	48.758	46.458	43.188	40.481	39.699	37.484	35.534
80	80.926	79.334	77.763	76.187	74.582	72.915	71.144	69.206	66.993	64.277	60.391	57.153	56.212	53.540	51.171
90	91.023	89.334	87.666	85.992	84.285	82.511	80.624	78.558	76.195	73.291	69.126	65.646	64.634	61.754	59.196
100	101.11	99.334	97.574	95.807	94.004	92.128	90.133	87.945	85.440	82.358	77.929	74.221	73.142	70.064	67.327
120	121.28	119.33	117.40	115.46	113.48	111.41	109.21	106.80	104.03	100.62	95.704	91.572	90.366	86.923	83.851
140	141.44	139.33	137.24	135.14	133.00	130.76	128.38	125.75	122.74	119.02	113.65	109.13	107.81	104.03	100.65
n/α	0.450	0.500	0.550	0.600	0.650	0.700	0.750	0.800	0.850	0.900	0.950	0.975	0.980	0.990	0.995

Distribución F de Fischer-Snedecor

Abscisas $F_{m,n,\alpha}$ que deixan á súa dereita un área α nunha F de Fischer-Snedecor con (m, n) graos de liberdade

$$\int_{F_{m,n,\alpha}}^{\infty} c_{m,n} x^{\frac{m}{2}-1} (n+mx)^{-\frac{m+n}{2}} dx = \alpha$$

$\alpha = 0,1$



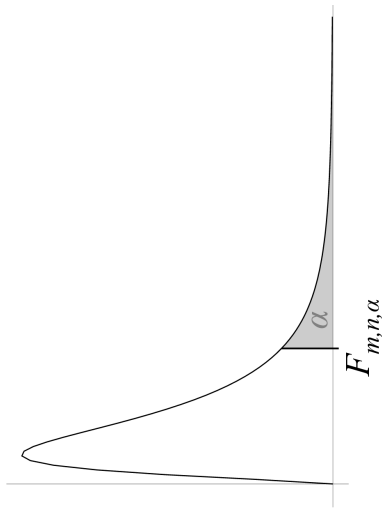
n\m	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	50	100	∞
1	39.863	49.500	53.593	55.833	57.240	58.204	58.906	59.439	59.858	60.195	60.705	61.220	61.740	62.002	62.265	62.688	63.007	63.328
2	8.526	9.000	9.162	9.243	9.293	9.326	9.349	9.367	9.381	9.392	9.408	9.425	9.441	9.450	9.458	9.471	9.481	9.491
3	5.538	5.462	5.391	5.343	5.309	5.285	5.266	5.252	5.240	5.230	5.216	5.200	5.184	5.176	5.168	5.155	5.144	5.134
4	4.545	4.325	4.191	4.107	4.051	4.010	3.979	3.955	3.936	3.920	3.896	3.870	3.844	3.831	3.817	3.795	3.778	3.761
5	4.060	3.780	3.619	3.520	3.453	3.405	3.368	3.339	3.316	3.297	3.268	3.238	3.207	3.191	3.174	3.147	3.126	3.105
6	3.776	3.463	3.289	3.181	3.108	3.055	3.014	2.983	2.958	2.937	2.905	2.871	2.836	2.818	2.800	2.770	2.746	2.722
7	3.589	3.257	3.074	2.961	2.883	2.827	2.785	2.752	2.725	2.703	2.668	2.632	2.595	2.575	2.555	2.523	2.497	2.471
8	3.458	3.113	2.924	2.806	2.726	2.668	2.624	2.589	2.561	2.538	2.502	2.464	2.425	2.404	2.383	2.348	2.321	2.293
9	3.360	3.006	2.813	2.693	2.611	2.551	2.505	2.469	2.440	2.416	2.379	2.340	2.298	2.277	2.255	2.218	2.189	2.159
10	3.285	2.924	2.728	2.605	2.522	2.461	2.414	2.377	2.347	2.323	2.284	2.244	2.201	2.178	2.155	2.117	2.087	2.055
11	3.225	2.860	2.660	2.536	2.451	2.389	2.342	2.304	2.274	2.248	2.209	2.167	2.123	2.100	2.076	2.036	2.005	1.972
12	3.177	2.807	2.606	2.480	2.394	2.331	2.283	2.245	2.214	2.188	2.147	2.105	2.060	2.036	2.011	1.970	1.938	1.904
13	3.136	2.763	2.560	2.434	2.347	2.283	2.234	2.195	2.164	2.138	2.097	2.053	2.007	1.983	1.958	1.915	1.882	1.846
14	3.102	2.726	2.522	2.395	2.307	2.243	2.193	2.154	2.122	2.095	2.054	2.010	1.962	1.938	1.912	1.869	1.834	1.797
15	3.073	2.695	2.490	2.361	2.273	2.208	2.158	2.119	2.086	2.059	2.017	1.972	1.924	1.899	1.873	1.828	1.793	1.755
16	3.048	2.668	2.462	2.333	2.244	2.178	2.128	2.088	2.055	2.028	1.985	1.940	1.891	1.866	1.839	1.793	1.757	1.718
17	3.026	2.645	2.437	2.308	2.218	2.152	2.102	2.061	2.028	2.001	1.958	1.912	1.862	1.836	1.809	1.763	1.726	1.686
18	3.007	2.624	2.416	2.286	2.196	2.130	2.079	2.038	2.005	1.977	1.933	1.887	1.837	1.810	1.783	1.736	1.698	1.657
19	2.990	2.606	2.397	2.266	2.176	2.109	2.058	2.017	1.984	1.956	1.912	1.865	1.814	1.787	1.759	1.711	1.673	1.631
20	2.975	2.589	2.380	2.249	2.158	2.091	2.040	1.999	1.965	1.937	1.892	1.845	1.794	1.767	1.738	1.690	1.650	1.607
21	2.961	2.575	2.365	2.233	2.142	2.075	2.023	1.982	1.948	1.920	1.875	1.827	1.776	1.748	1.719	1.670	1.630	1.586
22	2.949	2.561	2.351	2.219	2.128	2.060	2.008	1.967	1.933	1.904	1.859	1.811	1.759	1.731	1.702	1.652	1.611	1.567
23	2.937	2.549	2.339	2.207	2.115	2.047	1.995	1.953	1.919	1.890	1.845	1.796	1.744	1.716	1.686	1.636	1.594	1.549
24	2.927	2.538	2.327	2.195	2.103	2.035	1.983	1.941	1.906	1.877	1.832	1.783	1.730	1.702	1.672	1.621	1.579	1.533
25	2.918	2.528	2.317	2.184	2.092	2.024	1.971	1.929	1.895	1.866	1.820	1.771	1.718	1.689	1.659	1.607	1.565	1.518
26	2.909	2.519	2.307	2.174	2.082	2.014	1.961	1.919	1.884	1.855	1.809	1.760	1.706	1.677	1.647	1.594	1.551	1.504
27	2.901	2.511	2.299	2.165	2.073	2.005	1.952	1.909	1.874	1.845	1.799	1.749	1.695	1.666	1.636	1.583	1.539	1.491
28	2.894	2.503	2.291	2.157	2.064	1.996	1.943	1.900	1.865	1.836	1.790	1.740	1.685	1.656	1.625	1.572	1.528	1.478
29	2.887	2.495	2.283	2.149	2.057	1.988	1.935	1.892	1.857	1.827	1.781	1.731	1.676	1.647	1.616	1.562	1.517	1.467
30	2.881	2.489	2.276	2.142	2.049	1.980	1.927	1.884	1.849	1.819	1.773	1.722	1.667	1.638	1.606	1.552	1.507	1.456
40	2.835	2.440	2.226	2.091	1.997	1.927	1.873	1.829	1.793	1.763	1.715	1.662	1.605	1.574	1.541	1.483	1.434	1.377
60	2.791	2.393	2.177	2.041	1.946	1.875	1.819	1.775	1.738	1.707	1.657	1.603	1.543	1.511	1.476	1.413	1.358	1.291
100	2.756	2.356	2.139	2.002	1.906	1.834	1.778	1.732	1.695	1.663	1.612	1.557	1.494	1.460	1.423	1.355	1.293	1.214
200	2.731	2.329	2.111	1.973	1.876	1.804	1.747	1.701	1.663	1.631	1.579	1.522	1.458	1.422	1.383	1.310	1.242	1.144
∞	2.706	2.303	2.084	1.945	1.847	1.774	1.717	1.670	1.632	1.599	1.546	1.487	1.421	1.383	1.342	1.263	1.185	1.000
n/m	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	50	100	∞

Distribución F de Fischer-Snedecor

Abscisas $F_{m,n,\alpha}$ que deixan á súa dereita un área α nunha F de Fischer-Snedecor con (m, n) graos de liberdade

$$\int_{F_{m,n,\alpha}}^{\infty} c_{m,n} x^{\frac{m}{2}-1} (n+mx)^{-\frac{m+n}{2}} dx = \alpha$$

$\alpha = 0,05$



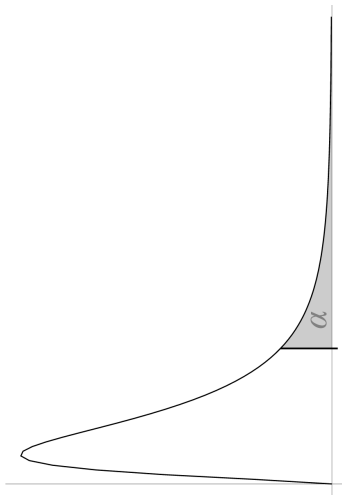
n \ m	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	50	100	∞
1	161.448	199.500	215.707	224.583	230.162	233.986	236.768	238.883	240.543	241.882	243.906	245.950	248.013	249.052	250.095	251.774	253.041	254.314
2	18.513	19.000	19.164	19.247	19.296	19.330	19.353	19.371	19.385	19.396	19.413	19.429	19.446	19.454	19.462	19.476	19.486	19.496
3	10.128	9.552	9.277	9.117	9.013	8.941	8.887	8.845	8.812	8.786	8.745	8.703	8.660	8.639	8.617	8.581	8.554	8.526
4	7.709	6.944	6.591	6.388	6.256	6.163	6.094	6.041	5.999	5.964	5.912	5.858	5.803	5.774	5.746	5.699	5.664	5.628
5	6.608	5.786	5.409	5.192	5.050	4.950	4.876	4.818	4.772	4.735	4.678	4.619	4.558	4.527	4.496	4.444	4.405	4.365
6	5.987	5.143	4.757	4.534	4.387	4.284	4.207	4.147	4.099	4.060	4.000	3.938	3.874	3.841	3.808	3.754	3.712	3.669
7	5.591	4.737	4.347	4.120	3.972	3.866	3.787	3.726	3.677	3.637	3.575	3.511	3.445	3.410	3.376	3.319	3.275	3.230
8	5.318	4.459	4.066	3.838	3.687	3.581	3.500	3.438	3.388	3.347	3.284	3.218	3.150	3.115	3.079	3.020	2.975	2.928
9	5.117	4.256	3.863	3.633	3.482	3.374	3.293	3.230	3.179	3.137	3.073	3.006	2.936	2.900	2.864	2.803	2.756	2.707
10	4.965	4.103	3.708	3.478	3.326	3.217	3.135	3.072	3.020	2.978	2.913	2.845	2.774	2.737	2.700	2.637	2.588	2.538
11	4.844	3.982	3.587	3.357	3.204	3.095	3.012	2.948	2.896	2.854	2.788	2.719	2.646	2.609	2.570	2.507	2.457	2.404
12	4.747	3.885	3.490	3.259	3.106	2.996	2.913	2.849	2.796	2.753	2.687	2.617	2.544	2.505	2.466	2.401	2.350	2.296
13	4.667	3.806	3.411	3.179	3.025	2.915	2.832	2.767	2.714	2.671	2.604	2.533	2.459	2.420	2.380	2.314	2.261	2.206
14	4.600	3.739	3.344	3.112	2.958	2.848	2.764	2.699	2.646	2.602	2.534	2.463	2.388	2.349	2.308	2.241	2.187	2.131
15	4.543	3.682	3.287	3.056	2.901	2.790	2.707	2.641	2.588	2.544	2.475	2.403	2.328	2.288	2.247	2.178	2.123	2.066
16	4.494	3.634	3.239	3.007	2.852	2.741	2.657	2.591	2.538	2.494	2.425	2.352	2.276	2.235	2.194	2.124	2.068	2.010
17	4.451	3.592	3.197	2.965	2.810	2.699	2.614	2.548	2.494	2.450	2.381	2.308	2.230	2.189	2.148	2.077	2.020	1.960
18	4.414	3.555	3.160	2.928	2.773	2.661	2.577	2.510	2.456	2.412	2.342	2.269	2.191	2.150	2.107	2.035	1.978	1.917
19	4.381	3.522	3.127	2.895	2.740	2.628	2.544	2.477	2.423	2.378	2.308	2.234	2.155	2.114	2.071	1.999	1.940	1.878
20	4.351	3.493	3.098	2.866	2.711	2.599	2.514	2.447	2.393	2.348	2.278	2.203	2.124	2.082	2.039	1.966	1.907	1.843
21	4.325	3.467	3.072	2.840	2.685	2.573	2.488	2.420	2.366	2.321	2.250	2.176	2.096	2.054	2.010	1.936	1.876	1.812
22	4.301	3.443	3.048	2.817	2.661	2.549	2.464	2.397	2.342	2.297	2.226	2.151	2.071	2.028	1.984	1.909	1.849	1.783
23	4.279	3.422	3.028	2.796	2.640	2.528	2.442	2.375	2.320	2.275	2.204	2.128	2.048	2.005	1.961	1.885	1.823	1.757
24	4.260	3.403	3.009	2.776	2.620	2.508	2.423	2.355	2.300	2.255	2.183	2.108	2.027	1.984	1.939	1.863	1.800	1.733
25	4.242	3.385	2.991	2.759	2.603	2.490	2.405	2.337	2.282	2.236	2.165	2.089	2.007	1.964	1.919	1.842	1.779	1.711
26	4.225	3.369	2.975	2.743	2.587	2.474	2.388	2.321	2.265	2.220	2.148	2.072	1.990	1.946	1.901	1.823	1.760	1.691
27	4.210	3.354	2.960	2.728	2.572	2.459	2.373	2.305	2.250	2.204	2.132	2.056	1.974	1.930	1.884	1.806	1.742	1.672
28	4.196	3.340	2.947	2.714	2.558	2.445	2.359	2.291	2.236	2.190	2.118	2.041	1.959	1.915	1.869	1.790	1.725	1.654
29	4.183	3.328	2.934	2.701	2.545	2.432	2.346	2.278	2.223	2.177	2.104	2.027	1.945	1.901	1.854	1.775	1.710	1.638
30	4.171	3.316	2.922	2.690	2.534	2.421	2.334	2.266	2.211	2.165	2.092	2.015	1.932	1.887	1.841	1.761	1.695	1.622
40	4.085	3.232	2.839	2.606	2.449	2.336	2.249	2.180	2.124	2.077	2.003	1.924	1.839	1.793	1.744	1.660	1.589	1.509
60	4.001	3.150	2.758	2.525	2.368	2.254	2.167	2.097	2.040	1.993	1.917	1.836	1.748	1.700	1.649	1.559	1.481	1.389
100	3.936	3.087	2.696	2.463	2.305	2.191	2.103	2.032	1.975	1.927	1.850	1.768	1.676	1.627	1.573	1.477	1.392	1.283
200	3.888	3.041	2.650	2.417	2.259	2.144	2.056	1.985	1.927	1.878	1.801	1.717	1.623	1.572	1.516	1.415	1.321	1.189
∞	3.841	2.996	2.605	2.372	2.214	2.099	2.010	1.938	1.880	1.831	1.752	1.666	1.571	1.517	1.459	1.350	1.243	1.000
n / m	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	50	100	∞

Distribución F de Fischer-Snedecor

Abscisas $F_{m,n,\alpha}$ que deixan á súa dereita un área α nunha F de Fisher-Snedecor con (m, n) graos de liberdade

$$\int_{F_{m,n,\alpha}}^{\infty} c_{m,n} x^{\frac{m}{2}-1} (n + mx)^{-\frac{m+n}{2}} dx = \alpha$$

$\alpha = 0,025$



$F_{m,n,\alpha}$

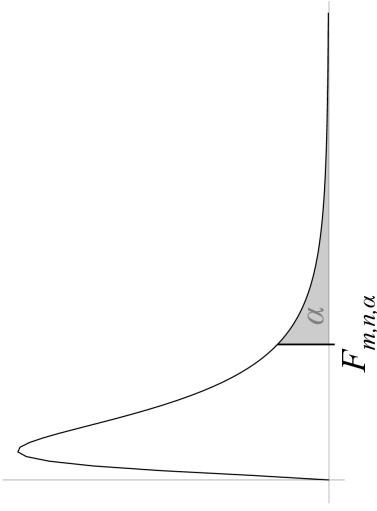
n \ m	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	50	100	∞
1	647.789	799.500	864.163	899.583	921.848	937.111	948.217	956.656	963.285	968.627	976.708	984.867	993.103	997.249	1001.414	1008.117	1013.175	1018.258
2	38.506	39.000	39.165	39.248	39.298	39.331	39.355	39.373	39.387	39.398	39.415	39.431	39.448	39.456	39.465	39.478	39.488	39.498
3	17.443	16.044	15.439	15.101	14.885	14.735	14.624	14.540	14.473	14.419	14.337	14.253	14.167	14.124	14.081	14.010	13.956	13.902
4	12.218	10.649	9.979	9.605	9.364	9.197	9.074	8.980	8.905	8.844	8.751	8.657	8.560	8.511	8.461	8.381	8.319	8.257
5	10.007	8.434	7.764	7.388	7.146	6.978	6.853	6.757	6.681	6.619	6.525	6.428	6.329	6.278	6.227	6.144	6.080	6.015
6	8.813	7.260	6.599	6.227	5.988	5.820	5.695	5.600	5.523	5.461	5.366	5.269	5.168	5.117	5.065	4.980	4.915	4.849
7	8.073	6.542	5.890	5.523	5.285	5.119	4.995	4.899	4.823	4.761	4.666	4.568	4.467	4.415	4.362	4.276	4.210	4.142
8	7.571	6.059	5.416	5.053	4.817	4.652	4.529	4.433	4.357	4.295	4.200	4.101	3.999	3.947	3.894	3.807	3.739	3.670
9	7.209	5.715	5.078	4.718	4.484	4.320	4.197	4.102	4.026	3.964	3.868	3.769	3.667	3.614	3.560	3.472	3.403	3.333
10	6.937	5.456	4.826	4.468	4.236	4.072	3.950	3.855	3.779	3.717	3.621	3.522	3.419	3.365	3.311	3.221	3.152	3.080
11	6.724	5.256	4.630	4.275	4.044	3.881	3.759	3.664	3.588	3.526	3.430	3.330	3.226	3.173	3.118	3.027	2.956	2.883
12	6.554	5.096	4.474	4.121	3.891	3.728	3.607	3.512	3.436	3.374	3.277	3.177	3.073	3.019	2.963	2.871	2.800	2.725
13	6.414	4.965	4.347	3.996	3.767	3.604	3.483	3.388	3.312	3.250	3.153	3.053	2.948	2.893	2.837	2.744	2.673	2.595
14	6.298	4.857	4.242	3.892	3.663	3.501	3.380	3.285	3.209	3.147	3.050	2.949	2.844	2.789	2.732	2.638	2.565	2.487
15	6.200	4.765	4.153	3.804	3.576	3.415	3.293	3.199	3.123	3.060	2.963	2.862	2.756	2.701	2.644	2.549	2.474	2.395
16	6.115	4.687	4.077	3.729	3.502	3.341	3.219	3.125	3.049	2.986	2.889	2.788	2.681	2.625	2.568	2.472	2.396	2.316
17	6.042	4.619	4.011	3.665	3.438	3.277	3.156	3.061	2.985	2.922	2.825	2.723	2.616	2.560	2.502	2.405	2.329	2.247
18	5.978	4.560	3.954	3.608	3.382	3.221	3.100	3.005	2.929	2.866	2.769	2.667	2.559	2.503	2.445	2.347	2.269	2.187
19	5.922	4.508	3.903	3.559	3.333	3.172	3.051	2.956	2.880	2.817	2.720	2.617	2.509	2.452	2.394	2.295	2.217	2.133
20	5.871	4.461	3.859	3.515	3.289	3.128	3.007	2.913	2.837	2.774	2.676	2.573	2.464	2.408	2.349	2.249	2.170	2.085
21	5.827	4.420	3.819	3.475	3.250	3.090	2.969	2.874	2.798	2.735	2.637	2.534	2.425	2.368	2.308	2.208	2.128	2.042
22	5.786	4.383	3.783	3.440	3.215	3.055	2.934	2.839	2.763	2.700	2.602	2.498	2.389	2.331	2.272	2.171	2.090	2.003
23	5.750	4.349	3.750	3.408	3.183	3.023	2.902	2.808	2.731	2.668	2.570	2.466	2.357	2.299	2.239	2.137	2.056	1.968
24	5.717	4.319	3.721	3.379	3.155	2.995	2.874	2.779	2.703	2.640	2.541	2.437	2.327	2.269	2.209	2.107	2.024	1.935
25	5.686	4.291	3.694	3.353	3.129	2.969	2.848	2.753	2.677	2.613	2.515	2.411	2.300	2.242	2.182	2.079	1.996	1.906
26	5.659	4.265	3.670	3.329	3.105	2.945	2.824	2.729	2.653	2.590	2.491	2.387	2.276	2.217	2.157	2.053	1.969	1.878
27	5.633	4.242	3.647	3.307	3.083	2.923	2.802	2.707	2.631	2.568	2.469	2.364	2.253	2.195	2.133	2.029	1.945	1.853
28	5.610	4.221	3.626	3.286	3.063	2.903	2.782	2.687	2.611	2.547	2.448	2.344	2.232	2.174	2.112	2.007	1.922	1.829
29	5.588	4.201	3.607	3.267	3.044	2.884	2.763	2.669	2.592	2.529	2.430	2.325	2.213	2.154	2.092	1.987	1.901	1.807
30	5.568	4.182	3.589	3.250	3.026	2.867	2.746	2.651	2.575	2.511	2.412	2.307	2.195	2.136	2.074	1.968	1.882	1.787
40	5.424	4.051	3.463	3.126	2.904	2.744	2.624	2.529	2.452	2.388	2.288	2.182	2.068	2.007	1.943	1.832	1.741	1.637
60	5.286	3.925	3.343	3.008	2.786	2.627	2.507	2.412	2.334	2.270	2.169	2.061	1.944	1.882	1.815	1.699	1.599	1.482
100	5.179	3.828	3.250	2.917	2.696	2.537	2.417	2.321	2.244	2.179	2.077	1.968	1.849	1.784	1.715	1.592	1.483	1.347
200	5.100	3.758	3.182	2.850	2.630	2.472	2.351	2.256	2.178	2.113	2.010	1.900	1.780	1.712	1.640	1.511	1.393	1.229
∞	5.024	3.689	3.116	2.786	2.567	2.408	2.288	2.192	2.114	2.048	1.945	1.833	1.708	1.640	1.566	1.428	1.296	1.000
n \ m	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	50	100	∞

Distribución F de Fischer-Snedecor

Abcissas $F_{m,n,\alpha}$ que deixan á súa dereita un área α nunha F de Fischer-Snedecor con (m, n) graos de liberdade

$$\int_{F_{m,n,\alpha}}^{\infty} c_{m,n} x^{\frac{m}{2}-1} (n+mx)^{-\frac{m+n}{2}} dx = \alpha$$

$\alpha = 0,01$



$F_{m,n,\alpha}$

n \ m	1	2	3	4	5	6	7	8	9	10	12	15	20	24	30	50	100	∞
1	4052.181	4999.500	5403.352	5624.583	5763.650	5858.986	5928.356	5981.070	6022.473	6055.847	6106.321	6157.285	6208.730	6234.631	6260.649	6302.517	6334.110	6365.864
2	98.503	99.000	99.166	99.249	99.299	99.333	99.356	99.374	99.388	99.399	99.416	99.433	99.449	99.458	99.466	99.479	99.489	99.499
3	34.116	30.817	29.457	28.710	28.237	27.911	27.672	27.489	27.345	27.229	27.052	26.872	26.690	26.598	26.505	26.354	26.240	26.125
4	21.198	18.000	16.694	15.977	15.522	15.207	14.976	14.799	14.659	14.546	14.374	14.198	14.020	13.929	13.838	13.690	13.577	13.463
5	16.258	13.274	12.060	11.392	10.967	10.672	10.456	10.289	10.158	10.051	9.888	9.722	9.553	9.466	9.379	9.238	9.130	9.020
6	13.745	10.925	9.780	9.148	8.746	8.466	8.260	8.102	7.976	7.874	7.718	7.559	7.396	7.313	7.229	7.091	6.987	6.880
7	12.246	9.547	8.451	7.847	7.460	7.191	6.993	6.840	6.719	6.620	6.469	6.314	6.155	6.074	5.992	5.858	5.755	5.650
8	11.259	8.649	7.591	7.006	6.632	6.371	6.178	6.029	5.911	5.814	5.667	5.515	5.359	5.279	5.198	5.065	4.963	4.859
9	10.561	8.022	6.992	6.422	6.057	5.802	5.613	5.467	5.351	5.257	5.111	4.962	4.808	4.729	4.649	4.517	4.415	4.311
10	10.044	7.559	6.552	5.994	5.636	5.386	5.200	5.057	4.942	4.849	4.706	4.558	4.405	4.327	4.247	4.115	4.014	3.909
11	9.646	7.206	6.217	5.668	5.316	5.069	4.886	4.744	4.632	4.539	4.397	4.251	4.099	4.021	3.941	3.810	3.708	3.602
12	9.330	6.927	5.953	5.412	5.064	4.821	4.640	4.499	4.388	4.296	4.155	4.010	3.858	3.780	3.701	3.569	3.467	3.361
13	9.074	6.701	5.739	5.205	4.862	4.620	4.441	4.302	4.191	4.100	3.960	3.815	3.665	3.587	3.507	3.375	3.272	3.165
14	8.862	6.515	5.564	5.035	4.695	4.456	4.278	4.140	4.030	3.939	3.800	3.656	3.505	3.427	3.348	3.215	3.112	3.004
15	8.683	6.359	5.417	4.893	4.556	4.318	4.142	4.004	3.895	3.805	3.666	3.522	3.372	3.294	3.214	3.081	2.977	2.868
16	8.531	6.226	5.292	4.773	4.437	4.202	4.026	3.890	3.780	3.691	3.553	3.409	3.259	3.181	3.101	2.967	2.863	2.753
17	8.400	6.112	5.185	4.669	4.336	4.102	3.927	3.791	3.682	3.593	3.455	3.312	3.162	3.084	3.003	2.869	2.764	2.653
18	8.285	6.013	5.092	4.579	4.248	4.015	3.841	3.705	3.597	3.508	3.371	3.227	3.077	2.999	2.919	2.784	2.678	2.566
19	8.185	5.926	5.010	4.500	4.171	3.939	3.765	3.631	3.523	3.434	3.297	3.153	3.003	2.925	2.844	2.709	2.602	2.489
20	8.096	5.849	4.938	4.431	4.103	3.871	3.699	3.564	3.457	3.368	3.231	3.088	2.938	2.859	2.778	2.643	2.535	2.421
21	8.017	5.780	4.874	4.369	4.042	3.812	3.640	3.506	3.398	3.310	3.173	3.030	2.880	2.801	2.720	2.584	2.475	2.360
22	7.945	5.719	4.817	4.313	3.988	3.758	3.587	3.453	3.346	3.258	3.121	2.978	2.827	2.749	2.667	2.531	2.422	2.305
23	7.881	5.664	4.765	4.264	3.939	3.710	3.539	3.406	3.299	3.211	3.074	2.931	2.781	2.702	2.620	2.483	2.373	2.256
24	7.823	5.614	4.718	4.218	3.895	3.667	3.496	3.363	3.256	3.168	3.032	2.889	2.738	2.659	2.577	2.440	2.329	2.211
25	7.770	5.568	4.675	4.177	3.855	3.627	3.457	3.324	3.217	3.129	2.993	2.850	2.699	2.620	2.538	2.400	2.289	2.169
26	7.721	5.526	4.637	4.140	3.818	3.591	3.421	3.288	3.182	3.094	2.958	2.815	2.664	2.585	2.503	2.364	2.252	2.131
27	7.677	5.488	4.601	4.106	3.785	3.558	3.388	3.256	3.149	3.062	2.926	2.783	2.632	2.552	2.470	2.330	2.218	2.097
28	7.636	5.453	4.568	4.074	3.754	3.528	3.358	3.226	3.120	3.032	2.896	2.753	2.602	2.522	2.440	2.300	2.187	2.064
29	7.598	5.420	4.538	4.045	3.725	3.499	3.330	3.198	3.092	3.005	2.868	2.726	2.574	2.495	2.412	2.271	2.158	2.034
30	7.562	5.390	4.510	4.018	3.699	3.473	3.304	3.173	3.067	2.979	2.843	2.700	2.549	2.469	2.386	2.245	2.131	2.006
40	7.314	5.179	4.313	3.828	3.514	3.291	3.124	2.993	2.888	2.801	2.665	2.522	2.369	2.288	2.203	2.058	1.938	1.805
60	7.077	4.977	4.126	3.649	3.339	3.119	2.953	2.823	2.718	2.632	2.496	2.352	2.198	2.115	2.028	1.877	1.749	1.601
100	6.895	4.824	3.984	3.513	3.206	2.988	2.823	2.694	2.590	2.503	2.368	2.223	2.067	1.983	1.893	1.735	1.598	1.427
200	6.763	4.713	3.881	3.414	3.110	2.893	2.730	2.601	2.497	2.411	2.275	2.129	1.971	1.886	1.794	1.629	1.481	1.279
∞	6.635	4.605	3.782	3.319	3.017	2.802	2.639	2.511	2.407	2.321	2.185	2.039	1.878	1.791	1.696	1.523	1.358	1.000

Bibliografía

- [1] J. S. Milton, *Estadística para Biología y Ciencias de la Salud* 3ª Edición. McGraw-Hill Interamericana, Madrid, 2007.
- [2] M. Samuels, J. A. Witmer, A. Schaffner, *Fundamentos de Estadística para las Ciencias de la Vida* 4ª Edición. Pearson Educación, S.A., Madrid, 2012.