

Notas sobre Regresión Logística
Breve introducción con aplicaciones

Carlos E Martinez-Rodriguez

Contents

Chapter 1

Principios

1.1 Introducción

La Estadística es una ciencia formal que estudia la recolección, análisis e interpretación de datos de una muestra representativa, ya sea para ayudar en la toma de decisiones o para explicar condiciones regulares o irregulares de algún fenómeno o estudio aplicado, de ocurrencia en forma aleatoria o condicional. Sin embargo, la estadística es más que eso, es decir, es transversal a una amplia variedad de disciplinas, desde la física hasta las ciencias sociales, desde las ciencias de la salud hasta el control de calidad. Se usa para la toma de decisiones en áreas de negocios o instituciones gubernamentales. Ahora bien, las técnicas estadísticas se aplican de manera amplia en mercadotecnia, contabilidad, control de calidad y en otras actividades; estudios de consumidores; análisis de resultados en deportes; administradores de instituciones; en la educación; organismos políticos; médicos; y por otras personas que intervienen en la toma de decisiones.

Definición 1. *La Estadística es la ciencia cuyo objetivo es reunir una información cuantitativa concerniente a individuos, grupos, series de hechos, etc. y deducir de ello gracias al análisis de estos datos unos significados precisos o unas previsiones para el futuro.*

La estadística, en general, es la ciencia que trata de la recopilación, organización presentación, análisis e interpretación de datos numéricos con el fin de realizar una toma de decisión más efectiva. Los métodos estadísticos tradicionalmente se utilizan para propósitos descriptivos, para organizar y resumir datos numéricos. La estadística descriptiva, por ejemplo trata de la tabulación de datos, su presentación en forma gráfica o ilustrativa y el cálculo de medidas descriptivas.

1.2 Historia de la Estadística

Es difícil conocer los orígenes de la Estadística. Desde los comienzos de la civilización han existido formas sencillas de estadística, pues ya se utilizaban representaciones gráficas y otros símbolos en pieles, rocas, palos de madera y paredes de cuevas para contar el número de personas, animales o ciertas cosas. Su origen empieza posiblemente en la isla de Cerdeña, donde existen monumentos prehistóricos pertenecientes a los Nuragas, los primeros habitantes de la isla; estos monumentos constan de bloques de basalto superpuestos sin mortero y en cuyas paredes se encontraban grabados toscos signos que han sido interpretados con mucha verosimilitud como muescas que servían para llevar la cuenta del ganado y la caza. Los babilonios usaban ya pequeñas tablillas de arcilla para recopilar datos en tablas sobre la producción agrícola y los géneros vendidos o cambiados mediante trueque. Otros vestigios pueden ser hallados en el antiguo Egipto, cuyos faraones lograron recopilar, hacia el año 3050 antes de

Cristo, prolijos datos relativos a la población y la riqueza del país. De acuerdo al historiador griego Heródoto, dicho registro de riqueza y población se hizo con el objetivo de preparar la construcción de las pirámides. En el mismo Egipto, Ramsés II hizo un censo de las tierras con el objeto de verificar un nuevo reparto. En el antiguo Israel la Biblia da referencias, en el libro de los Números, de los datos estadísticos obtenidos en dos recuentos de la población hebrea. El rey David por otra parte, ordenó a Joab, general del ejército hacer un censo de Israel con la finalidad de conocer el número de la población. También los chinos efectuaron censos hace más de cuarenta siglos. Los griegos efectuaron censos periódicamente con fines tributarios, sociales (división de tierras) y militares (cálculo de recursos y hombres disponibles). La investigación histórica revela que se realizaron 69 censos para calcular los impuestos, determinar los derechos de voto y ponderar la potencia guerrera.

Fueron los romanos, maestros de la organización política, quienes mejor supieron emplear los recursos de la estadística. Cada cinco años realizaban un censo de la población y sus funcionarios públicos tenían la obligación de anotar nacimientos, defunciones y matrimonios, sin olvidar los recuentos periódicos del ganado y de las riquezas contenidas en las tierras conquistadas. Para el nacimiento de Cristo sucedía uno de estos empadronamientos de la población bajo la autoridad del imperio. Durante los mil años siguientes a la caída del imperio Romano se realizaron muy pocas operaciones Estadísticas, con la notable excepción de las relaciones de tierras pertenecientes a la Iglesia, compiladas por Pipino el Breve en el 758 y por Carlomagno en el 762 DC. Durante el siglo IX se realizaron en Francia algunos censos parciales de siervos. En Inglaterra, Guillermo el Conquistador recopiló el Domesday Book o libro del Gran Catastro para el año 1086, un documento de la propiedad, extensión y valor de las tierras de Inglaterra. Esa obra fue el primer compendio estadístico de Inglaterra. Aunque Carlomagno, en Francia; y Guillermo el Conquistador, en Inglaterra, trataron de revivir la técnica romana, los métodos estadísticos permanecieron casi olvidados durante la Edad Media.

Durante los siglos XV, XVI, y XVII, hombres como Leonardo de Vinci, Nicolás Copérnico, Galileo, Neper, William Harvey, Sir Francis Bacon y René Descartes, hicieron grandes operaciones al método científico, de tal forma que cuando se crearon los Estados Nacionales y surgió como fuerza el comercio internacional existía ya un método capaz de aplicarse a los datos económicos. Para el año 1532 empezaron a registrarse en Inglaterra las defunciones debido al temor que Enrique VII tenía por la peste. Más o menos por la misma época, en Francia la ley exigió a los clérigos registrar los bautismos, fallecimientos y matrimonios. Durante un brote de peste que apareció a fines de la década de 1500, el gobierno inglés comenzó a publicar estadística semanales de los decesos. Esa costumbre continuó muchos años, y en 1632 estos Bills of Mortality (Cuentas de Mortalidad) contenían los nacimientos y fallecimientos por sexo. En 1662, el capitán John Graunt usó documentos que abarcaban treinta años y efectuó predicciones sobre el número de personas que morirían de varias enfermedades y sobre las proporciones de nacimientos de varones y mujeres que cabría esperar. El trabajo de Graunt, condensado en su obra *Natural and Political Observations...Made upon the Bills of Mortality*, fue un esfuerzo innovador en el análisis estadístico. Por el año 1540 el alemán Sebastián Muster realizó una compilación estadística de los recursos nacionales, comprensiva de datos sobre organización política, instrucciones sociales, comercio y poderío militar.

Los eruditos del siglo XVII demostraron especial interés por la Estadística Demográfica como resultado de la especulación sobre si la población aumentaba, decrecía o permanecía estática. En los tiempos modernos tales métodos fueron resucitados por algunos reyes que necesitaban conocer las riquezas monetarias y el potencial humano de sus respectivos países. El primer empleo de los datos estadísticos para fines ajenos a la política tuvo lugar en 1691 y estuvo a cargo de Gaspar Neumann, un profesor alemán que vivía en Breslau. Este investigador se propuso destruir la antigua creencia popular de que en los años terminados en siete moría más gente que en los restantes, y para lograrlo hurgó pacientemente en los archivos parroquiales de la ciudad. Después de revisar miles de partidas de defunción pudo demostrar que en tales años no fallecían más personas que en los demás. Los procedimientos de Neumann

fueron conocidos por el astrónomo inglés Halley, descubridor del cometa que lleva su nombre, quien los aplicó al estudio de la vida humana.

Durante el siglo XVII y principios del XVIII, matemáticos como Bernoulli, Francis Maseres, Laplace y Laplace desarrollaron la teoría de probabilidades. No obstante durante cierto tiempo, la teoría de las probabilidades limitó su aplicación a los juegos de azar y hasta el siglo XVIII no comenzó a aplicarse a los grandes problemas científicos. Godofredo Achenwall, profesor de la Universidad de Gotinga, acuñó en 1760 la palabra estadística, que extrajo del término italiano *statista* (estadista). Creía, y con sobrada razón, que los datos de la nueva ciencia serían el aliado más eficaz del gobernante consciente. La raíz remota de la palabra se halla, por otra parte, en el término latino *status*, que significa estado o situación; Esta etimología aumenta el valor intrínseco de la palabra, por cuanto la estadística revela el sentido cuantitativo de las más variadas situaciones. Jacques Quételet es quien aplica las Estadísticas a las ciencias sociales. Este interpretó la teoría de la probabilidad para su uso en las ciencias sociales y resolver la aplicación del principio de promedios y de la variabilidad a los fenómenos sociales. Quételet fue el primero en realizar la aplicación práctica de todo el método Estadístico, entonces conocido, a las diversas ramas de la ciencia. Entretanto, en el período del 1800 al 1820 se desarrollaron dos conceptos matemáticos fundamentales para la teoría Estadística; la teoría de los errores de observación, aportada por Laplace y Gauss; y la teoría de los mínimos cuadrados desarrollada por Laplace, Gauss y Legendre. A finales del siglo XIX, Sir Francis Gaston ideó el método conocido por Correlación, que tenía por objeto medir la influencia relativa de los factores sobre las variables. De aquí partió el desarrollo del coeficiente de correlación creado por Karl Pearson y otros cultivadores de la ciencia biométrica como J. Pease Norton, R. H. Hooker y G. Udny Yule, que efectuaron amplios estudios sobre la medida de las relaciones.

La historia de la estadística está resumida en tres grandes etapas o fases.

- **Fase 1: Los Censos:** Desde el momento en que se constituye una autoridad política, la idea de inventariar de una forma más o menos regular la población y las riquezas existentes en el territorio está ligada a la conciencia de soberanía y a los primeros esfuerzos administrativos.
- **Fase 2: De la Descripción de los Conjuntos a la Aritmética Política:** Las ideas mercantilistas extrañan una intensificación de este tipo de investigación. Colbert multiplica las encuestas sobre artículos manufacturados, el comercio y la población: los intendentes del Reino envían a París sus memorias. Vauban, más conocido por sus fortificaciones o su *Dime Royale*, que es la primera propuesta de un impuesto sobre los ingresos, se señala como el verdadero precursor de los sondeos. Más tarde, Bufón se preocupa de esos problemas antes de dedicarse a la historia natural. La escuela inglesa proporciona un nuevo progreso al superar la fase puramente descriptiva.

Sus tres principales representantes son Graunt, Petty y Halley. El penúltimo es autor de la famosa *Aritmética Política*. Chaptal, ministro del interior francés, publica en 1801 el primer censo general de población, desarrolla los estudios industriales, de las producciones y los cambios, haciéndose sistemáticos durante las dos terceras partes del siglo XIX.

- **Fase 3: Estadística y Cálculo de Probabilidades:** El cálculo de probabilidades se incorpora rápidamente como un instrumento de análisis extremadamente poderoso para el estudio de los fenómenos económicos y sociales y en general para el estudio de fenómenos cuyas causas son demasiadas complejas para conocerlos totalmente y hacer posible su análisis.

La Estadística para su mejor estudio se ha dividido en dos grandes ramas: **la Estadística Descriptiva y la Estadística Inferencial**.

- **Descriptiva:** consiste sobre todo en la presentación de datos en forma de tablas y gráficas. Esta comprende cualquier actividad relacionada con los datos y está diseñada para resumir o describir

los mismos sin factores pertinentes adicionales; esto es, sin intentar inferir nada que vaya más allá de los datos, como tales.

- **Inferencial:** se deriva de muestras, de observaciones hechas sólo acerca de una parte de un conjunto numeroso de elementos y esto implica que su análisis requiere de generalizaciones que van más allá de los datos. Como consecuencia, la característica más importante del reciente crecimiento de la estadística ha sido un cambio en el énfasis de los métodos que describen a métodos que sirven para hacer generalizaciones. La Estadística Inferencial investiga o analiza una población partiendo de una muestra tomada.

Estadística Inferencial

Los métodos básicos de la estadística inferencial son la estimación y el contraste de hipótesis, que juegan un papel fundamental en la investigación. Por tanto, algunos de los objetivos que se persiguen son:

- Calcular los parámetros de la distribución de medias o proporciones muestrales de tamaño n , extraídas de una población de media y varianza conocidas.
- Estimar la media o la proporción de una población a partir de la media o proporción muestral.
- Utilizar distintos tamaños muestrales para controlar la confianza y el error admitido.
- Contrastar los resultados obtenidos a partir de muestras.
- Visualizar gráficamente, mediante las respectivas curvas normales, las estimaciones realizadas.

En definitiva, la idea es, a partir de una población se extrae una muestra por algunos de los métodos existentes, con la que se generan datos numéricos que se van a utilizar para generar estadísticos con los que realizar estimaciones o contrastes poblacionales. Existen dos formas de estimar parámetros: la *estimación puntual* y la *estimación por intervalo de confianza*. En la primera se busca, con base en los datos muestrales, un único valor estimado para el parámetro. Para la segunda, se determina un intervalo dentro del cual se encuentra el valor del parámetro, con una probabilidad determinada.

Si el objetivo del tratamiento estadístico inferencial, es efectuar generalizaciones acerca de la estructura, composición o comportamiento de las poblaciones no observadas, a partir de una parte de la población, será necesario que la proporción de población examinada sea representativa del total. Por ello, la selección de la muestra requiere unos requisitos que lo garanticen, debe ser representativa y aleatoria.

Además, la cantidad de elementos que integran la muestra (el tamaño de la muestra) depende de múltiples factores, como el dinero y el tiempo disponibles para el estudio, la importancia del tema analizado, la confiabilidad que se espera de los resultados, las características propias del fenómeno analizado, etcétera.

Así, a partir de la muestra seleccionada se realizan algunos cálculos y se estima el valor de los parámetros de la población tales como la media, la varianza, la desviación estándar, o la forma de la distribución, etc.

El conjunto de los métodos que se utilizan para medir las características de la información, para resumir los valores individuales, y para analizar los datos a fin de extraerles el máximo de información, es lo que se llama *métodos estadísticos*. Los métodos de análisis para la información cuantitativa se pueden dividir en los siguientes seis pasos:

- Definición del problema.

- Recopilación de la información existente.
- Obtención de información original.
- Clasificación.
- Presentación.
- Análisis.

El centro de gravedad de la metodología estadística se empieza a desplazar técnicas de computación intensiva aplicadas a grandes masas de datos, y se empieza a considerar el método estadístico como un proceso iterativo de búsqueda del modelo ideal. Las aplicaciones en este periodo de la Estadística a la Economía conducen a una disciplina con contenido propio: la Econometría. La investigación estadística en problemas militares durante la segunda guerra mundial y los nuevos métodos de programación matemática, dan lugar a la Investigación Operativa. El tratamiento de los datos de la investigación científica tiene varias etapas:

- En la etapa de recolección de datos del método científico, se define a la población de interés y se selecciona una muestra o conjunto de personas representativas de la misma, se realizan experimentos o se emplean instrumentos ya existentes o de nueva creación, para medir los atributos de interés necesarios para responder a las preguntas de investigación. Durante lo que es llamado trabajo de campo se obtienen los datos en crudo, es decir las respuestas directas de los sujetos uno por uno, se codifican (se les asignan valores a las respuestas), se capturan y se verifican para ser utilizados en las siguientes etapas.
- En la etapa de recuento, se organizan y ordenan los datos obtenidos de la muestra. Esta será descrita en la siguiente etapa utilizando la estadística descriptiva, todas las investigaciones utilizan estadística descriptiva, para conocer de manera organizada y resumida las características de la muestra.
- En la etapa de análisis se utilizan las pruebas estadísticas (estadística inferencial) y en la interpretación se acepta o rechaza la hipótesis nula.

Niveles de medición y tipos de variables

Para poder emplear el método estadístico en un estudio es necesario medir las variables.

- Medir: es asignar valores a las propiedades de los objetos bajo ciertas reglas, esas reglas son los niveles de medición.
- Cuantificar: es asignar valores a algo tomando un patrón de referencia. Por ejemplo, cuantificar es ver cuántos hombres y cuántas mujeres hay.

Variable: es una característica o propiedad que asume diferentes valores dentro de una población de interés y cuya variación es susceptible de medirse.

Las variables pueden clasificarse de acuerdo al tipo de valores que puede tomar como:

- **Discretas o categóricas** en las que los valores se relacionan a nombres, etiquetas o categorías, no existe un significado numérico directo.

- **Continuas** los valores tienen un correlato numérico directo, son continuos y susceptibles de fraccionarse y de poder utilizarse en operaciones aritméticas.
- **Dicotómica** sólo tienen dos valores posibles, la característica está ausente o presente.

En cuanto a una clasificación estadística, las variables pueden ser:

- **Aleatoria** Aquella en la cual desconocemos el valor porque fluctúa de acuerdo a un evento debido al azar.
- **Determinística** Aquella variable de la que se conoce el valor.
- **Independiente** aquellas variables que son manipuladas por el investigador. Define los grupos.
- **Dependiente** son mediciones que ocurren durante el experimento o tratamiento (resultado de la independiente), es la que se mide y compara entre los grupos.

En lo que tiene que ver con los **Niveles de Medición** tenemos distintos tipos de variable

- **Nominal:** Las propiedades de la medición nominal son:
 - Exhaustiva: implica a todas las opciones.
 - A los sujetos se les asignan categorías, por lo que son mutuamente excluyentes. Es decir, la variable está presente o no; tiene o no una característica.
- **Ordinal:** Las propiedades de la medición ordinal son:
 - El nivel ordinal posee transitividad, por lo que se tiene la capacidad de identificar que es mejor o mayor que otra, en ese sentido se pueden establecer jerarquías.
 - Las distancias entre un valor y otro no son iguales.
- **Intervalo:**
 - El nivel de medición intervalar requiere distancias iguales entre cada valor. Por lo general utiliza datos cuantitativos. Por ejemplo: temperatura, atributos psicológicos (CI, nivel de autoestima, pruebas de conocimientos, etc.)
 - Las unidades de calificación son equivalentes en todos los puntos de la escala. Una escala de intervalos implica: clasificación, magnitud y unidades de tamaños iguales (Brown, 2000).
 - Se pueden hacer operaciones aritméticas.
 - Cuando se le pide al sujeto que califique una situación del 0 al 10 puede tomarse como un nivel de medición de intervalo, siempre y cuando se incluya el 0.
- **Razón:**
 - La escala empieza a partir del 0 absoluto, por lo tanto incluye sólo los números por su valor en sí, por lo que no pueden existir los números con signo negativo. Por ejemplo: Peso corporal en kg., edad en años, estatura en cm.

Definiciones adicionales

- **Variable:** Consideraciones que una variable son una característica o fenómeno que puede tomar distintos valores.
- **Dato:** Mediciones o cualidades que han sido recopiladas como resultado de observaciones.
- **Población:** Se considera el área de la cual son extraídos los datos. Es decir, es el conjunto de elementos o individuos que poseen una característica común y medible acerca de lo cual se desea información. Es también llamado Universo.
- **Muestra:** Es un subconjunto de la población, seleccionado de acuerdo a una regla o algún plan de muestreo.
- **Censo:** Recopilación de todos los datos (de interés para la investigación) de la población.
- **Estadística:** Es una función o fórmula que depende de los datos de la muestra (es variable).
- **Parámetro:** Característica medible de la población. Es un resumen numérico de alguna variable observada de la población. Los parámetros normales que se estudian son: *La media poblacional, Proporción.*
- **Estimador:** Un estimador de un parámetro es un estadístico que se emplea para conocer el parámetro desconocido.
- **Estadístico:** Es una función de los valores de la muestra. Es una variable aleatoria, cuyos valores dependen de la muestra seleccionada. Su distribución de probabilidad, se conoce como *Distribución muestral del estadístico.*
- **Estimación:** Este término indica que a partir de lo observado en una muestra (un resumen estadístico con las medidas que conocemos de Descriptiva) se extrapola o generaliza dicho resultado muestral a la población total, de modo que lo estimado es el valor generalizado a la población. Consiste en la búsqueda del valor de los parámetros poblacionales objeto de estudio. Puede ser puntual o por intervalo de confianza:
 - *Puntual:* cuando buscamos un valor concreto. Un estimador de un parámetro poblacional es una función de los datos muestrales. En pocas palabras, es una fórmula que depende de los valores obtenidos de una muestra, para realizar estimaciones. Lo que se pretende obtener es el valor exacto de un parámetro.
 - *Intervalo de confianza:* cuando determinamos un intervalo, dentro del cual se supone que va a estar el valor del parámetro que se busca con una cierta probabilidad. El intervalo de confianza está determinado por dos valores dentro de los cuales afirmamos que está el verdadero parámetro con cierta probabilidad. Son unos límites o margen de variabilidad que damos al valor estimado, para poder afirmar, bajo un criterio de probabilidad, que el verdadero valor no los rebasará.

Este intervalo contiene al parámetro estimado con una determinada certeza o nivel de confianza.

En la estimación por intervalos se usan los siguientes conceptos:

- **Variabilidad del parámetro:** Si no se conoce, puede obtenerse una aproximación en los datos o en un estudio piloto. También hay métodos para calcular el tamaño de la muestra que prescinden de este aspecto. Habitualmente se usa como medida de esta variabilidad la desviación típica poblacional.
- **Error de la estimación:** Es una medida de su precisión que se corresponde con la amplitud del intervalo de confianza. Cuanta más precisión se desee en la estimación de un parámetro, más estrecho deberá ser el intervalo de confianza y, por tanto, menor el error, y más sujetos deberán incluirse en la muestra estudiada.
- **Nivel de confianza:** Es la probabilidad de que el verdadero valor del parámetro estimado en la población se sitúe en el intervalo de confianza obtenido. El nivel de confianza se denota por $1 - \alpha$
- **p -value :** También llamado nivel de significación. Es la probabilidad (en tanto por uno) de fallar en nuestra estimación, esto es, la diferencia entre la certeza (1) y el nivel de confianza $1 - \alpha$.
- **Valor crítico:** Se representa por $Z_{\alpha/2}$. Es el valor de la abscisa en una determinada distribución que deja a su derecha un área igual a $1/2$, siendo $1 - \alpha$ el nivel de confianza. Normalmente los valores críticos están tabulados o pueden calcularse en función de la distribución de la población.

Para un tamaño fijo de la muestra, los conceptos de error y nivel de confianza van relacionados. Si admitimos un error mayor, esto es, aumentamos el tamaño del intervalo de confianza, tenemos también una mayor probabilidad de éxito en nuestra estimación, es decir, un mayor nivel de confianza. Por tanto, un aspecto que debe de tenerse en cuenta es el tamaño muestral, ya que para disminuir el error que se comente habrá que aumentar el tamaño muestral. Esto se resolverá, para un intervalo de confianza cualquiera, despejando el tamaño de la muestra en cualquiera de las formulas de los intervalos de confianza que veremos a continuación, a partir del error máximo permitido. Los intervalos de confianza pueden ser unilaterales o bilaterales:

- **Contraste de Hipótesis:** Consiste en determinar si es aceptable, partiendo de datos muestrales, que la característica o el parámetro poblacional estudiado tome un determinado valor o esté dentro de unos determinados valores.
- **Nivel de Confianza:** Indica la proporción de veces que acertaríamos al afirmar que el parámetro está dentro del intervalo al seleccionar muchas muestras.

1.3 Muestreo:

Muestreo: Una muestra es representativa en la medida que es imagen de la población. En general, podemos decir que el tamaño de una muestra dependerá principalmente de: *Nivel de precisión deseado, Recursos disponibles, Tiempo involucrado en la investigación*. Además el plan de muestreo debe considerar *La población, Parámetros a medir*. Existe una gran cantidad de tipos de muestreo, en la práctica los más utilizados son los siguientes:

- **MUESTREO ALEATORIO SIMPLE:** Es un método de selección de n unidades extraídas de N , de tal manera que cada una de las posibles muestras tiene la misma probabilidad de ser escogida. (En la práctica, se enumeran las unidades de 1 a N , y a continuación se seleccionan n números aleatorios entre 1 y N , ya sea de tablas o de alguna urna con fichas numeradas).

- **MUESTREO ESTRATIFICADO ALEATORIO:** Se usa cuando la población está agrupada en pocos estratos, cada uno de ellos son muchas entidades. Este muestreo consiste en sacar una muestra aleatoria simple de cada uno de los estratos. (Generalmente, de tamaño proporcional al estrato).
- **MUESTREO SISTEMÁTICO:** Se utiliza cuando las unidades de la población están de alguna manera totalmente ordenadas. Para seleccionar una muestra de n unidades, se divide la población en n subpoblaciones de tamaño $K = N/n$ y se toma al azar una unidad de la K primeras y de ahí en adelante cada K -ésima unidad.
- **MUESTREO POR CONGLOMERADO:** Se emplea cuando la población está dividida en grupos o conglomerados pequeños. Consiste en obtener una muestra aleatoria simple de conglomerados y luego CENSAR cada uno de éstos.
- **MUESTREO EN DOS ETAPAS (Bietápico):** En este caso la muestra se toma en dos pasos:
 - Seleccionar una muestra de unidades primarias, y
 - Seleccionar una muestra de elementos a partir de cada unidad primaria escogida.
 - *Observación:* En la realidad es posible encontrarse con situaciones en las cuales no es posible aplicar libremente un tipo de muestreo, incluso estaremos obligados a mezclarlas en ocasiones.

1.4 Errores Estadísticos Comunes

El propósito de esta sección es solamente indicar los malos usos comunes de datos estadísticos, sin incluir el uso de métodos estadísticos complicados. Un estudiante debería estar alerta en relación con estos malos usos y debería hacer un gran esfuerzo para evitarlos a fin de ser un verdadero estadístico.

Datos estadísticos inadecuados: Los datos estadísticos son usados como la materia prima para un estudio estadístico. Cuando los datos son inadecuados, la conclusión extraída del estudio de los datos se vuelve obviamente inválida. Por ejemplo, supongamos que deseamos encontrar el ingreso familiar típico del año pasado en la ciudad Y de 50,000 familias y tenemos una muestra consistente del ingreso de solamente tres familias: 1 millón, 2 millones y no ingreso. Si sumamos el ingreso de las tres familias y dividimos el total por 3, obtenemos un promedio de 1 millón. Entonces, extraemos una conclusión basada en la muestra de que el ingreso familiar promedio durante el año pasado en la ciudad fue de 1 millón. Es obvio que la conclusión es falsa, puesto que las cifras son extremas y el tamaño de la muestra es demasiado pequeño; por lo tanto la muestra no es representativa.

Hay muchas otras clases de datos inadecuados. Por ejemplo, algunos datos son respuestas inexactas de una encuesta, porque las preguntas usadas en la misma son vagas o engañosas, algunos datos son toscas estimaciones porque no hay disponibles datos exactos o es demasiado costosa su obtención, y algunos datos son irrelevantes en un problema dado, porque el estudio estadístico no está bien planeado. Al momento de recopilar los datos que serán procesados se es susceptible de cometer errores así como durante los cálculos de los mismos. No obstante, hay otros errores que no tienen nada que ver con la digitación y que no son tan fácilmente identificables. Algunos de éstos errores son:

- **Sesgo:** Es imposible ser completamente objetivo o no tener ideas preconcebidas antes de comenzar a estudiar un problema, y existen muchas maneras en que una perspectiva o estado mental pueda influir en la recopilación y en el análisis de la información. En estos casos se dice que hay un sesgo cuando el individuo da mayor peso a los datos que apoyan su opinión que a aquellos que la contradicen. Un caso extremo de sesgo sería la situación donde primero se toma una decisión y después se utiliza el análisis estadístico para justificar la decisión ya tomada.

- **Datos No Comparables:** el establecer comparaciones es una de las partes más importantes del análisis estadístico, pero es extremadamente importante que tales comparaciones se hagan entre datos que sean comparables.
- **Proyección descuidada de tendencias:** la proyección simplista de tendencias pasadas hacia el futuro es uno de los errores que más ha desacreditado el uso del análisis estadístico.
- **Muestreo Incorrecto:** en la mayoría de los estudios sucede que el volumen de información disponible es tan inmenso que se hace necesario estudiar muestras, para derivar conclusiones acerca de la población a que pertenece la muestra. Si la muestra se selecciona correctamente, tendrá básicamente las mismas propiedades que la población de la cual fue extraída; pero si el muestreo se realiza incorrectamente, entonces puede suceder que los resultados no signifiquen nada.

Sesgo significa que un usuario dé los datos perjudicialmente de más énfasis a los hechos, los cuales son empleados para mantener su predeterminada posición u opinión. Los estadísticos son frecuentemente degradados por lemas tales como: *Hay tres clases de mentiras: mentiras, mentiras reprobables y estadística, y Las cifras no mienten, pero los mentirosos piensan.* Hay dos clases de sesgos: conscientes e inconscientes. Ambos son comunes en el análisis estadístico. Hay numerosos ejemplos de sesgos conscientes. Un anunciante frecuentemente usa la estadística para probar que su producto es muy superior al producto de su competidor. Un político prefiere usar la estadística para sostener su punto de vista. Gerentes y líderes de trabajadores pueden simultáneamente situar sus respectivas cifras estadísticas sobre la misma tabla de trato para mostrar que sus rechazos o peticiones son justificadas. Es casi imposible que un sesgo inconsciente esté completamente ausente en un trabajo estadístico. En lo que respecta al ser humano, es difícil obtener una actitud completamente objetiva al abordar un problema, aun cuando un científico debería tener una mente abierta. Un estadístico debería estar enterado del hecho de que su interpretación de los resultados del análisis estadístico está influenciado por su propia experiencia, conocimiento y antecedentes con relación al problema dado.

Chapter 2

documento: Bases

2.1 Análisis de Regresion Lineal (RL)

Nota 1. • En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.
- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

2.2 Análisis de Regresion Lineal (RL)

Nota 2. • En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.
- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

2.2.1 Regresión Lineal Simple (RLS)

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (2.1)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (2.2)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación (??) se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (2.3)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (2.4)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (2.5)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (2.6)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (2.7)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (2.8)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (2.9)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (2.10)$$

2.2.2 Regresión Lineal Simple (RLS)

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (2.11)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (2.12)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación (??) se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (2.13)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (2.14)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales

resultante:

$$\begin{aligned}\frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0\end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned}-2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0\end{aligned}$$

simplificando

$$\begin{aligned}n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i\end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (2.15)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (2.16)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (2.17)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (2.18)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (2.19)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (2.20)$$

2.3 3. Análisis de Regresión Lineal (RL)

Nota 3. • En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.

2.3.1 3.1 Regresión Lineal Simple (RLS)

- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (2.21)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

2.3.2 3.2 Método de Mínimos Cuadrados

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (2.22)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación ?? se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (2.23)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (2.24)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (2.25)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (2.26)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (2.27)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (2.28)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (2.29)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (2.30)$$

2.3.3 3.3 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Nota 4. • Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo.

- Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$E(\hat{\beta}_1) = E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right)$$

$$\begin{aligned}
&= \frac{1}{S_{xx}} E \left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x}) \right) \\
&= \frac{1}{S_{xx}} \left[\beta_0 E \left(\sum_{k=1}^n (x_k - \bar{x}) \right) + E \left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x}) \right) \right. \\
&\quad \left. + E \left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x}) \right) \right] = \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1
\end{aligned}$$

por lo tanto

$$E(\hat{\beta}_1) = \beta_1 \quad (2.31)$$

Nota 5. Es decir, $\hat{\beta}_1$ es un estimador insesgado.

Ahora calculemos la varianza:

$$\begin{aligned}
V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) \\
&= \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) = \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 \\
&= \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}}
\end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad (2.32)$$

Proposición 1.

$$\begin{aligned}
E(\hat{\beta}_0) &= \beta_0, \\
V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \\
Cov(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{\sigma^2 \bar{x}}{S_{xx}}.
\end{aligned}$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$$e_k = y_k - \hat{y}_k, \text{ se le denomina } \mathbf{residuo}.$$

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \quad (2.33)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$\begin{aligned}
SC_E &= \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \\
E(SC_E) &= (n-2)\sigma^2, \text{ por lo tanto} \\
\hat{\sigma}^2 &= \frac{SC_E}{n-2} = MC_E \text{ es un estimador insesgado de } \sigma^2.
\end{aligned}$$

2.3.4 3.4 Prueba de Hipótesis en RLS

- Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza.
- Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$.

Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_1, \sigma^2)$.

De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??). Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (2.34)$$

que se distribuye t con $n - 2$ grados de libertad bajo $H_0 : \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (2.35)$$

Para β_0 se puede proceder de manera análoga para

$$H_0: \beta_0 = \beta_{0,0},$$

$$H_1: \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}, \quad (2.36)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (2.37)$$

- No rechazar $H_0 : \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0 : \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0 : \beta_1 = 0$ puede realizarse de la siguiente manera:

$$\begin{aligned}
S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\
S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 \\
&= \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\
&= \sum_{k=1}^n [(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2] \\
&= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\
&= \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) \\
&= \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\
&= \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&= \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&= \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0.
\end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (2.38)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (2.39)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (2.40)$$

Por lo tanto la ecuación (??) se puede reescribir como

$$S_{yy} = SC_R + SC_E \quad (2.41)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$

$$\begin{aligned} S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}) \\ S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R \end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 2.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx} \quad (2.42)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E}$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$.

El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n - 2$	MC_E	
Total	S_{yy}	$n - 1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}} \quad (2.43)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

2.3.5 Estimación de Intervalos en RLS

- Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros.
- El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión.

Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n - 2$ grados de libertad. Por tanto un intervalo de confianza de $100(1 - \alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}}. \quad (2.44)$$

De igual manera, para β_0 un intervalo de confianza al $100(1 - \alpha)\%$ es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (2.45)$$

2.4 3. Análisis de Regresión Lineal (RL)

Nota 6. • En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.

2.4.1 3.1 Regresión Lineal Simple (RLS)

- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (2.46)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

2.4.2 3.2 Método de Mínimos Cuadrados

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (2.47)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación ?? se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (2.48)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (2.49)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (2.50)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (2.51)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (2.52)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (2.53)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (2.54)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (2.55)$$

2.4.3 3.3 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Nota 7. • Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo.

- Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$\begin{aligned}
 E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right) \\
 &= \frac{1}{S_{xx}} E\left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x})\right) \\
 &= \frac{1}{S_{xx}} \left[\beta_0 E\left(\sum_{k=1}^n (x_k - \bar{x})\right) + E\left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x})\right) \right. \\
 &\quad \left. + E\left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x})\right) \right] = \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1
 \end{aligned}$$

por lo tanto

$$E(\hat{\beta}_1) = \beta_1 \quad (2.56)$$

Nota 8. Es decir, $\hat{\beta}_1$ es un estimador insesgado.

Ahora calculemos la varianza:

$$\begin{aligned}
 V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) \\
 &= \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) = \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 \\
 &= \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}}
 \end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad (2.57)$$

Proposición 3.

$$\begin{aligned}
 E(\hat{\beta}_0) &= \beta_0, \\
 V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \\
 Cov(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{\sigma^2 \bar{x}}{S_{xx}}.
 \end{aligned}$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$e_k = y_k - \hat{y}_k$, se le denomina **residuo**.

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \quad (2.58)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$\begin{aligned} SC_E &= \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \\ E(SC_E) &= (n-2)\sigma^2, \text{ por lo tanto} \\ \hat{\sigma}^2 &= \frac{SC_E}{n-2} = MC_E \text{ es un estimador insesgado de } \sigma^2. \end{aligned}$$

2.4.4 3.4 Prueba de Hipótesis en RLS

- Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza.
- Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$.

Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_i, \sigma^2)$.

De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??).

Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \beta_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (2.59)$$

que se distribuye t con $n-2$ grados de libertad bajo $H_0: \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (2.60)$$

Para β_0 se puede proceder de manera análoga para

$$H_0: \beta_0 = \beta_{0,0},$$

$$H_1: \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}, \quad (2.61)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (2.62)$$

- No rechazar $H_0 : \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0 : \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0 : \beta_1 = 0$ puede realizarse de la siguiente manera:

$$\begin{aligned}
S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\
S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 \\
&= \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\
&= \sum_{k=1}^n \left[(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2 \right] \\
&= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\
&= \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) \\
&= \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\
&= \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&= \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&= \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0.
\end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (2.63)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (2.64)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (2.65)$$

Por lo tanto la ecuación (??) se puede reescribir como

$$S_{yy} = SC_R + SC_E \quad (2.66)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$

$$\begin{aligned} S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}) \\ S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R \end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 4.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx} \quad (2.67)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E}$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$.

El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n - 2$	MC_E	
Total	S_{yy}	$n - 1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}} \quad (2.68)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

2.4.5 Estimación de Intervalos en RLS

- Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros.
- El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión.

Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n - 2$ grados de libertad. Por tanto un intervalo de confianza de $100(1 - \alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}}. \quad (2.69)$$

De igual manera, para β_0 un intervalo de confianza al 100(1 - α) % es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (2.70)$$

2.4.6 Predicción

Supongamos que se tiene un valor x_0 de interés, entonces la estimación puntual de este nuevo valor

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 \quad (2.71)$$

Nota 9. Esta nueva observación es independiente de las utilizadas para obtener el modelo de regresión, por tanto, el intervalo en torno a la recta de regresión es inapropiado, puesto que se basa únicamente en los datos empleados para ajustar el modelo de regresión.

El intervalo de confianza en torno a la recta de regresión se refiere a la respuesta media verdadera $x = x_0$, no a observaciones futuras.

Sea y_0 la observación futura en $x = x_0$, y sea \hat{y}_0 dada en la ecuación anterior, el estimador de y_0 . Si se define la variable aleatoria

$$w = y_0 - \hat{y}_0,$$

esta se distribuye normalmente con media cero y varianza

$$V(w) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]$$

dado que y_0 es independiente de \hat{y}_0 , por lo tanto el intervalo de predicción al nivel α para futuras observaciones x_0 es

$$\begin{aligned} \hat{y}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]} &\leq y_0 \\ &\leq \hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]}. \end{aligned}$$

2.4.7 Prueba de falta de ajuste

Es común encontrar que el modelo ajustado no satisface totalmente el modelo necesario para los datos, en este caso es preciso saber qué tan bueno es el modelo propuesto. Para esto se propone la siguiente prueba de hipótesis:

H_0 : El modelo propuesto se ajusta adecuadamente a los datos.

H_1 : El modelo NO se ajusta a los datos.

La prueba implica dividir la suma de cuadrados del error o del residuo en las siguientes dos componentes:

$$SC_E = SC_{EP} + SC_{FDA}$$

donde SC_{EP} es la suma de cuadrados atribuibles al error puro, y SC_{FDA} es la suma de cuadrados atribuible a la falta de ajuste del modelo.

2.4.8 Coeficiente de Determinación

La cantidad

$$R^2 = \frac{SC_R}{S_{yy}} = 1 - \frac{SC_E}{S_{yy}} \quad (2.72)$$

se denomina coeficiente de determinación y se utiliza para saber si el modelo de regresión es suficiente o no. Se puede demostrar que $0 \leq R^2 \leq 1$, una manera de interpretar este valor es que si $R^2 = k$, entonces el modelo de regresión explica el $k * 100\%$ de la variabilidad en los datos. R^2

- No mide la magnitud de la pendiente de la recta de regresión
- Un valor grande de R^2 no implica una pendiente empinada.
- No mide la suficiencia del modelo.
- Valores grandes de R^2 no implican necesariamente que el modelo de regresión proporcionará predicciones precisas para futuras observaciones.

Chapter 3

Introducción

3.1 Introducción

La Estadística es una ciencia formal que estudia la recolección, análisis e interpretación de datos de una muestra representativa, ya sea para ayudar en la toma de decisiones o para explicar condiciones regulares o irregulares de algún fenómeno o estudio aplicado, de ocurrencia en forma aleatoria o condicional. Sin embargo, la estadística es más que eso, es decir, es transversal a una amplia variedad de disciplinas, desde la física hasta las ciencias sociales, desde las ciencias de la salud hasta el control de calidad. Se usa para la toma de decisiones en áreas de negocios o instituciones gubernamentales. Ahora bien, las técnicas estadísticas se aplican de manera amplia en mercadotecnia, contabilidad, control de calidad y en otras actividades; estudios de consumidores; análisis de resultados en deportes; administradores de instituciones; en la educación; organismos políticos; médicos; y por otras personas que intervienen en la toma de decisiones.

Definición 2. *La Estadística es la ciencia cuyo objetivo es reunir una información cuantitativa concerniente a individuos, grupos, series de hechos, etc. y deducir de ello gracias al análisis de estos datos unos significados precisos o unas previsiones para el futuro.*

La estadística, en general, es la ciencia que trata de la recopilación, organización presentación, análisis e interpretación de datos numéricos con el fin de realizar una toma de decisión más efectiva. Los métodos estadísticos tradicionalmente se utilizan para propósitos descriptivos, para organizar y resumir datos numéricos. La estadística descriptiva, por ejemplo trata de la tabulación de datos, su presentación en forma gráfica o ilustrativa y el cálculo de medidas descriptivas.

3.1.1 Historia de la Estadística

Es difícil conocer los orígenes de la Estadística. Desde los comienzos de la civilización han existido formas sencillas de estadística, pues ya se utilizaban representaciones gráficas y otros símbolos en pieles, rocas, palos de madera y paredes de cuevas para contar el número de personas, animales o ciertas cosas. Su origen empieza posiblemente en la isla de Cerdeña, donde existen monumentos prehistóricos pertenecientes a los Nuragas, los primeros habitantes de la isla; estos monumentos constan de bloques de basalto superpuestos sin mortero y en cuyas paredes de encontraban grabados toscos signos que han sido interpretados con mucha verosimilitud como muescas que servían para llevar la cuenta del ganado y la caza. Los babilonios usaban ya pequeñas tablillas de arcilla para recopilar datos en tablas sobre la producción agrícola y los géneros vendidos o cambiados mediante trueque. Otros vestigios pueden ser hallados en el antiguo Egipto, cuyos faraones lograron recopilar, hacia el año 3050 antes de Cristo, prolijos datos relativos a la población y la riqueza del país. De acuerdo al historiador griego Heródoto, dicho registro de riqueza y población se hizo con el objetivo de preparar la construcción de las pirámides. En el mismo Egipto, Ramsés II hizo un censo de las tierras con el objeto de verificar un nuevo reparto. En el antiguo Israel la Biblia da referencias, en el libro de los Números, de los datos estadísticos obtenidos en dos recuentos de la población hebrea. El rey David por otra parte, ordenó a Joab,

general del ejército hacer un censo de Israel con la finalidad de conocer el número de la población. También los chinos efectuaron censos hace más de cuarenta siglos. Los griegos efectuaron censos periódicamente con fines tributarios, sociales (división de tierras) y militares (cálculo de recursos y hombres disponibles). La investigación histórica revela que se realizaron 69 censos para calcular los impuestos, determinar los derechos de voto y ponderar la potencia guerrera.

Fueron los romanos, maestros de la organización política, quienes mejor supieron emplear los recursos de la estadística. Cada cinco años realizaban un censo de la población y sus funcionarios públicos tenían la obligación de anotar nacimientos, defunciones y matrimonios, sin olvidar los recuentos periódicos del ganado y de las riquezas contenidas en las tierras conquistadas. Para el nacimiento de Cristo sucedía uno de estos empadronamientos de la población bajo la autoridad del imperio. Durante los mil años siguientes a la caída del imperio Romano se realizaron muy pocas operaciones Estadísticas, con la notable excepción de las relaciones de tierras pertenecientes a la Iglesia, compiladas por Pipino el Breve en el 758 y por Carlomagno en el 762 DC. Durante el siglo IX se realizaron en Francia algunos censos parciales de siervos. En Inglaterra, Guillermo el Conquistador recopiló el Domesday Book o libro del Gran Catastro para el año 1086, un documento de la propiedad, extensión y valor de las tierras de Inglaterra. Esa obra fue el primer compendio estadístico de Inglaterra. Aunque Carlomagno, en Francia; y Guillermo el Conquistador, en Inglaterra, trataron de revivir la técnica romana, los métodos estadísticos permanecieron casi olvidados durante la Edad Media.

Durante los siglos XV, XVI, y XVII, hombres como Leonardo de Vinci, Nicolás Copérnico, Galileo, Neper, William Harvey, Sir Francis Bacon y René Descartes, hicieron grandes operaciones al método científico, de tal forma que cuando se crearon los Estados Nacionales y surgió como fuerza el comercio internacional existía ya un método capaz de aplicarse a los datos económicos. Para el año 1532 empezaron a registrarse en Inglaterra las defunciones debido al temor que Enrique VII tenía por la peste. Más o menos por la misma época, en Francia la ley exigió a los clérigos registrar los bautismos, fallecimientos y matrimonios. Durante un brote de peste que apareció a fines de la década de 1500, el gobierno inglés comenzó a publicar estadística semanales de los decesos. Esa costumbre continuó muchos años, y en 1632 estos Bills of Mortality (Cuentas de Mortalidad) contenían los nacimientos y fallecimientos por sexo. En 1662, el capitán John Graunt usó documentos que abarcaban treinta años y efectuó predicciones sobre el número de personas que morirían de varias enfermedades y sobre las proporciones de nacimientos de varones y mujeres que cabría esperar. El trabajo de Graunt, condensado en su obra *Natural and Political Observations...Made upon the Bills of Mortality*, fue un esfuerzo innovador en el análisis estadístico. Por el año 1540 el alemán Sebastián Muster realizó una compilación estadística de los recursos nacionales, comprensiva de datos sobre organización política, instrucciones sociales, comercio y poderío militar.

Los eruditos del siglo XVII demostraron especial interés por la Estadística Demográfica como resultado de la especulación sobre si la población aumentaba, decrecía o permanecía estática. En los tiempos modernos tales métodos fueron resucitados por algunos reyes que necesitaban conocer las riquezas monetarias y el potencial humano de sus respectivos países. El primer empleo de los datos estadísticos para fines ajenos a la política tuvo lugar en 1691 y estuvo a cargo de Gaspar Neumann, un profesor alemán que vivía en Breslau. Este investigador se propuso destruir la antigua creencia popular de que en los años terminados en siete moría más gente que en los restantes, y para lograrlo hurgó pacientemente en los archivos parroquiales de la ciudad. Después de revisar miles de partidas de defunción pudo demostrar que en tales años no fallecían más personas que en los demás. Los procedimientos de Neumann fueron conocidos por el astrónomo inglés Halley, descubridor del cometa que lleva su nombre, quien los aplicó al estudio de la vida humana.

Durante el siglo XVII y principios del XVIII, matemáticos como Bernoulli, Francis Maseres, Lagrange y Laplace desarrollaron la teoría de probabilidades. No obstante durante cierto tiempo, la teoría de las probabilidades limitó su aplicación a los juegos de azar y hasta el siglo XVIII no comenzó a aplicarse a los grandes problemas científicos. Godofredo Achenwall, profesor de la Universidad de Gotinga, acuñó en 1760 la palabra estadística, que extrajo del término italiano statista (estadista). Creía, y con sobrada razón, que los datos de la nueva ciencia serían el aliado más eficaz del gobernante consciente. La raíz remota de la palabra se halla, por otra parte, en el término latino status, que significa estado o situación; Esta etimología aumenta el valor intrínseco de la palabra, por cuanto la estadística revela el sentido cuantitativo de las más variadas situaciones. Jacques Quételet es quien aplica las Estadísticas a las ciencias sociales. Este interpretó la teoría

de la probabilidad para su uso en las ciencias sociales y resolver la aplicación del principio de promedios y de la variabilidad a los fenómenos sociales. Quételet fue el primero en realizar la aplicación práctica de todo el método Estadístico, entonces conocido, a las diversas ramas de la ciencia. Entretanto, en el período del 1800 al 1820 se desarrollaron dos conceptos matemáticos fundamentales para la teoría Estadística; la teoría de los errores de observación, aportada por Laplace y Gauss; y la teoría de los mínimos cuadrados desarrollada por Laplace, Gauss y Legendre. A finales del siglo XIX, Sir Francis Gaston ideó el método conocido por Correlación, que tenía por objeto medir la influencia relativa de los factores sobre las variables. De aquí partió el desarrollo del coeficiente de correlación creado por Karl Pearson y otros cultivadores de la ciencia biométrica como J. Pease Norton, R. H. Hooker y G. Udny Yule, que efectuaron amplios estudios sobre la medida de las relaciones.

La historia de la estadística está resumida en tres grandes etapas o fases.

- **Fase 1: Los Censos:** Desde el momento en que se constituye una autoridad política, la idea de inventariar de una forma más o menos regular la población y las riquezas existentes en el territorio está ligada a la conciencia de soberanía y a los primeros esfuerzos administrativos.
- **Fase 2: De la Descripción de los Conjuntos a la Aritmética Política:** Las ideas mercantilistas extrañan una intensificación de este tipo de investigación. Colbert multiplica las encuestas sobre artículos manufacturados, el comercio y la población: los intendentes del Reino envían a París sus memorias. Vauban, más conocido por sus fortificaciones o su Dime Royale, que es la primera propuesta de un impuesto sobre los ingresos, se señala como el verdadero precursor de los sondeos. Más tarde, Bufón se preocupa de esos problemas antes de dedicarse a la historia natural. La escuela inglesa proporciona un nuevo progreso al superar la fase puramente descriptiva.

Sus tres principales representantes son Graunt, Petty y Halley. El penúltimo es autor de la famosa Aritmética Política. Chaptal, ministro del interior francés, publica en 1801 el primer censo general de población, desarrolla los estudios industriales, de las producciones y los cambios, haciéndose sistemáticos durante las dos terceras partes del siglo XIX.

- **Fase 3: Estadística y Cálculo de Probabilidades:** El cálculo de probabilidades se incorpora rápidamente como un instrumento de análisis extremadamente poderoso para el estudio de los fenómenos económicos y sociales y en general para el estudio de fenómenos cuyas causas son demasiado complejas para conocerlos totalmente y hacer posible su análisis.

La Estadística para su mejor estudio se ha dividido en dos grandes ramas: **la Estadística Descriptiva y la Estadística Inferencial.**

- **Descriptiva:** consiste sobre todo en la presentación de datos en forma de tablas y gráficas. Esta comprende cualquier actividad relacionada con los datos y está diseñada para resumir o describir los mismos sin factores pertinentes adicionales; esto es, sin intentar inferir nada que vaya más allá de los datos, como tales.
- **Inferencial:** se deriva de muestras, de observaciones hechas sólo acerca de una parte de un conjunto numeroso de elementos y esto implica que su análisis requiere de generalizaciones que van más allá de los datos. Como consecuencia, la característica más importante del reciente crecimiento de la estadística ha sido un cambio en el énfasis de los métodos que describen a métodos que sirven para hacer generalizaciones. La Estadística Inferencial investiga o analiza una población partiendo de una muestra tomada.

Estadística Inferencial

Los métodos básicos de la estadística inferencial son la estimación y el contraste de hipótesis, que juegan un papel fundamental en la investigación. Por tanto, algunos de los objetivos que se persiguen son:

- Calcular los parámetros de la distribución de medias o proporciones muestrales de tamaño n , extraídas de una población de media y varianza conocidas.
- Estimar la media o la proporción de una población a partir de la media o proporción muestral.
- Utilizar distintos tamaños muestrales para controlar la confianza y el error admitido.
- Contrastar los resultados obtenidos a partir de muestras.
- Visualizar gráficamente, mediante las respectivas curvas normales, las estimaciones realizadas.

En definitiva, la idea es, a partir de una población se extrae una muestra por algunos de los métodos existentes, con la que se generan datos numéricos que se van a utilizar para generar estadísticos con los que realizar estimaciones o contrastes poblacionales. Existen dos formas de estimar parámetros: la *estimación puntual* y la *estimación por intervalo de confianza*. En la primera se busca, con base en los datos muestrales, un único valor estimado para el parámetro. Para la segunda, se determina un intervalo dentro del cual se encuentra el valor del parámetro, con una probabilidad determinada.

Si el objetivo del tratamiento estadístico inferencial, es efectuar generalizaciones acerca de la estructura, composición o comportamiento de las poblaciones no observadas, a partir de una parte de la población, será necesario que la proporción de población examinada sea representativa del total. Por ello, la selección de la muestra requiere unos requisitos que lo garanticen, debe ser representativa y aleatoria.

Además, la cantidad de elementos que integran la muestra (el tamaño de la muestra) depende de múltiples factores, como el dinero y el tiempo disponibles para el estudio, la importancia del tema analizado, la confiabilidad que se espera de los resultados, las características propias del fenómeno analizado, etcétera.

Así, a partir de la muestra seleccionada se realizan algunos cálculos y se estima el valor de los parámetros de la población tales como la media, la varianza, la desviación estándar, o la forma de la distribución, etc.

El conjunto de los métodos que se utilizan para medir las características de la información, para resumir los valores individuales, y para analizar los datos a fin de extraerles el máximo de información, es lo que se llama *métodos estadísticos*. Los métodos de análisis para la información cuantitativa se pueden dividir en los siguientes seis pasos:

- Definición del problema.
- Recopilación de la información existente.
- Obtención de información original.
- Clasificación.
- Presentación.
- Análisis.

El centro de gravedad de la metodología estadística se empieza a desplazar técnicas de computación intensiva aplicadas a grandes masas de datos, y se empieza a considerar el método estadístico como un proceso iterativo de búsqueda del modelo ideal. Las aplicaciones en este periodo de la Estadística a la Economía conducen a una disciplina con contenido propio: la Econometría. La investigación estadística en problemas militares durante la segunda guerra mundial y los nuevos métodos de programación matemática, dan lugar a la Investigación Operativa. El tratamiento de los datos de la investigación científica tiene varias etapas:

- En la etapa de recolección de datos del método científico, se define a la población de interés y se selecciona una muestra o conjunto de personas representativas de la misma, se realizan experimentos o se emplean instrumentos ya existentes o de nueva creación, para medir los atributos de interés necesarios para responder a las preguntas de investigación. Durante lo que es llamado trabajo de campo se obtienen los datos en crudo, es decir las respuestas directas de los sujetos uno por uno, se codifican (se les asignan valores a las respuestas), se capturan y se verifican para ser utilizados en las siguientes etapas.

- En la etapa de recuento, se organizan y ordenan los datos obtenidos de la muestra. Esta será descrita en la siguiente etapa utilizando la estadística descriptiva, todas las investigaciones utilizan estadística descriptiva, para conocer de manera organizada y resumida las características de la muestra.
- En la etapa de análisis se utilizan las pruebas estadísticas (estadística inferencial) y en la interpretación se acepta o rechaza la hipótesis nula.

Niveles de medición y tipos de variables

Para poder emplear el método estadístico en un estudio es necesario medir las variables.

- **Medir:** es asignar valores a las propiedades de los objetos bajo ciertas reglas, esas reglas son los niveles de medición.
- **Cuantificar:** es asignar valores a algo tomando un patrón de referencia. Por ejemplo, cuantificar es ver cuántos hombres y cuántas mujeres hay.

Variable: es una característica o propiedad que asume diferentes valores dentro de una población de interés y cuya variación es susceptible de medirse.

Las variables pueden clasificarse de acuerdo al tipo de valores que puede tomar como:

- **Discretas o categóricas** en las que los valores se relacionan a nombres, etiquetas o categorías, no existe un significado numérico directo.
- **Continuas** los valores tienen un correlato numérico directo, son continuos y susceptibles de fraccionarse y de poder utilizarse en operaciones aritméticas.
- **Dicotómica** sólo tienen dos valores posibles, la característica está ausente o presente.

En cuanto a una clasificación estadística, las variables pueden ser:

- **Aleatoria** Aquella en la cual desconocemos el valor porque fluctúa de acuerdo a un evento debido al azar.
- **Determinística** Aquella variable de la que se conoce el valor.
- **Independiente** aquellas variables que son manipuladas por el investigador. Define los grupos.
- **Dependiente** son mediciones que ocurren durante el experimento o tratamiento (resultado de la independiente), es la que se mide y compara entre los grupos.

En lo que tiene que ver con los **Niveles de Medición** tenemos distintos tipos de variable

- **Nominal:** Las propiedades de la medición nominal son:
 - Exhaustiva: implica a todas las opciones.
 - A los sujetos se les asignan categorías, por lo que son mutuamente excluyentes. Es decir, la variable está presente o no; tiene o no una característica.
- **Ordinal:** Las propiedades de la medición ordinal son:
 - El nivel ordinal posee transitividad, por lo que se tiene la capacidad de identificar que es mejor o mayor que otra, en ese sentido se pueden establecer jerarquías.
 - Las distancias entre un valor y otro no son iguales.
- **Intervalo:**

- El nivel de medición intervalar requiere distancias iguales entre cada valor. Por lo general utiliza datos cuantitativos. Por ejemplo: temperatura, atributos psicológicos (CI, nivel de autoestima, pruebas de conocimientos, etc.)
- Las unidades de calificación son equivalentes en todos los puntos de la escala. Una escala de intervalos implica: clasificación, magnitud y unidades de tamaños iguales (Brown, 2000).
- Se pueden hacer operaciones aritméticas.
- Cuando se le pide al sujeto que califique una situación del 0 al 10 puede tomarse como un nivel de medición de intervalo, siempre y cuando se incluya el 0.

- **Razón:**

- La escala empieza a partir del 0 absoluto, por lo tanto incluye sólo los números por su valor en sí, por lo que no pueden existir los números con signo negativo. Por ejemplo: Peso corporal en kg., edad en años, estatura en cm.

Definiciones adicionales

- **Variable:** Consideraciones que una variable son una característica o fenómeno que puede tomar distintos valores.
- **Dato:** Mediciones o cualidades que han sido recopiladas como resultado de observaciones.
- **Población:** Se considera el área de la cual son extraídos los datos. Es decir, es el conjunto de elementos o individuos que poseen una característica común y medible acerca de lo cual se desea información. Es también llamado Universo.
- **Muestra:** Es un subconjunto de la población, seleccionado de acuerdo a una regla o algún plan de muestreo.
- **Censo:** Recopilación de todos los datos (de interés para la investigación) de la población.
- **Estadística:** Es una función o fórmula que depende de los datos de la muestra (es variable).
- **Parámetro:** Característica medible de la población. Es un resumen numérico de alguna variable observada de la población. Los parámetros normales que se estudian son: *La media poblacional, Proporción.*
- **Estimador:** Un estimador de un parámetro es un estadístico que se emplea para conocer el parámetro desconocido.
- **Estadístico:** Es una función de los valores de la muestra. Es una variable aleatoria, cuyos valores dependen de la muestra seleccionada. Su distribución de probabilidad, se conoce como *Distribución muestral del estadístico.*
- **Estimación:** Este término indica que a partir de lo observado en una muestra (un resumen estadístico con las medidas que conocemos de Descriptiva) se extrapola o generaliza dicho resultado muestral a la población total, de modo que lo estimado es el valor generalizado a la población. Consiste en la búsqueda del valor de los parámetros poblacionales objeto de estudio. Puede ser puntual o por intervalo de confianza:
 - **Puntual:** cuando buscamos un valor concreto. Un estimador de un parámetro poblacional es una función de los datos muestrales. En pocas palabras, es una fórmula que depende de los valores obtenidos de una muestra, para realizar estimaciones. Lo que se pretende obtener es el valor exacto de un parámetro.

- **Intervalo de confianza:** cuando determinamos un intervalo, dentro del cual se supone que va a estar el valor del parámetro que se busca con una cierta probabilidad. El intervalo de confianza está determinado por dos valores dentro de los cuales afirmamos que está el verdadero parámetro con cierta probabilidad. Son unos límites o margen de variabilidad que damos al valor estimado, para poder afirmar, bajo un criterio de probabilidad, que el verdadero valor no los rebasará.

Este intervalo contiene al parámetro estimado con una determinada certeza o nivel de confianza.

En la estimación por intervalos se usan los siguientes conceptos:

- **Variabilidad del parámetro:** Si no se conoce, puede obtenerse una aproximación en los datos o en un estudio piloto. También hay métodos para calcular el tamaño de la muestra que prescinden de este aspecto. Habitualmente se usa como medida de esta variabilidad la desviación típica poblacional.
- **Error de la estimación:** Es una medida de su precisión que se corresponde con la amplitud del intervalo de confianza. Cuanta más precisión se desee en la estimación de un parámetro, más estrecho deberá ser el intervalo de confianza y, por tanto, menor el error, y más sujetos deberán incluirse en la muestra estudiada.
- **Nivel de confianza:** Es la probabilidad de que el verdadero valor del parámetro estimado en la población se sitúe en el intervalo de confianza obtenido. El nivel de confianza se denota por $1 - \alpha$
- **p-value:** También llamado nivel de significación. Es la probabilidad (en tanto por uno) de fallar en nuestra estimación, esto es, la diferencia entre la certeza (1) y el nivel de confianza $1 - \alpha$.
- **Valor crítico:** Se representa por $Z_{\alpha/2}$. Es el valor de la abscisa en una determinada distribución que deja a su derecha un área igual a $1/2$, siendo $1 - \alpha$ el nivel de confianza. Normalmente los valores críticos están tabulados o pueden calcularse en función de la distribución de la población.

Para un tamaño fijo de la muestra, los conceptos de error y nivel de confianza van relacionados. Si admitimos un error mayor, esto es, aumentamos el tamaño del intervalo de confianza, tenemos también una mayor probabilidad de éxito en nuestra estimación, es decir, un mayor nivel de confianza. Por tanto, un aspecto que debe de tenerse en cuenta es el tamaño muestral, ya que para disminuir el error que se comente habrá que aumentar el tamaño muestral. Esto se resolverá, para un intervalo de confianza cualquiera, despejando el tamaño de la muestra en cualquiera de las formulas de los intervalos de confianza que veremos a continuación, a partir del error máximo permitido. Los intervalos de confianza pueden ser unilaterales o bilaterales:

- **Contraste de Hipótesis:** Consiste en determinar si es aceptable, partiendo de datos muestrales, que la característica o el parámetro poblacional estudiado tome un determinado valor o esté dentro de unos determinados valores.
- **Nivel de Confianza:** Indica la proporción de veces que acertaríamos al afirmar que el parámetro está dentro del intervalo al seleccionar muchas muestras.

3.1.2 Muestreo:

Muestreo: Una muestra es representativa en la medida que es imagen de la población. En general, podemos decir que el tamaño de una muestra dependerá principalmente de: *Nivel de precisión deseado, Recursos disponibles, Tiempo involucrado en la investigación*. Además el plan de muestreo debe considerar *La población, Parámetros a medir*. Existe una gran cantidad de tipos de muestreo, en la práctica los más utilizados son los siguientes:

- **MUESTREO ALEATORIO SIMPLE:** Es un método de selección de n unidades extraídas de N , de tal manera que cada una de las posibles muestras tiene la misma probabilidad de ser escogida. (En la práctica, se enumeran las unidades de 1 a N , y a continuación se seleccionan n números aleatorios entre 1 y N , ya sea de tablas o de alguna urna con fichas numeradas).

- **MUESTREO ESTRATIFICADO ALEATORIO:** Se usa cuando la población está agrupada en pocos estratos, cada uno de ellos son muchas entidades. Este muestreo consiste en sacar una muestra aleatoria simple de cada uno de los estratos. (Generalmente, de tamaño proporcional al estrato).
- **MUESTREO SISTEMÁTICO:** Se utiliza cuando las unidades de la población están de alguna manera totalmente ordenadas. Para seleccionar una muestra de n unidades, se divide la población en n subpoblaciones de tamaño $K = N/n$ y se toma al azar una unidad de la K primeras y de ahí en adelante cada K -ésima unidad.
- **MUESTREO POR CONGLOMERADO:** Se emplea cuando la población está dividida en grupos o conglomerados pequeños. Consiste en obtener una muestra aleatoria simple de conglomerados y luego CENSAR cada uno de éstos.
- **MUESTREO EN DOS ETAPAS (Bietápico):** En este caso la muestra se toma en dos pasos:
 - Seleccionar una muestra de unidades primarias, y
 - Seleccionar una muestra de elementos a partir de cada unidad primaria escogida.
 - *Observación:* En la realidad es posible encontrarse con situaciones en las cuales no es posible aplicar libremente un tipo de muestreo, incluso estaremos obligados a mezclarlas en ocasiones.

3.1.3 Errores Estadísticos Comunes

El propósito de esta sección es solamente indicar los malos usos comunes de datos estadísticos, sin incluir el uso de métodos estadísticos complicados. Un estudiante debería estar alerta en relación con estos malos usos y debería hacer un gran esfuerzo para evitarlos a fin de ser un verdadero estadístico.

Datos estadísticos inadecuados: Los datos estadísticos son usados como la materia prima para un estudio estadístico. Cuando los datos son inadecuados, la conclusión extraída del estudio de los datos se vuelve obviamente inválida. Por ejemplo, supongamos que deseamos encontrar el ingreso familiar típico del año pasado en la ciudad Y de 50,000 familias y tenemos una muestra consistente del ingreso de solamente tres familias: 1 millón, 2 millones y no ingreso. Si sumamos el ingreso de las tres familias y dividimos el total por 3, obtenemos un promedio de 1 millón. Entonces, extraemos una conclusión basada en la muestra de que el ingreso familiar promedio durante el año pasado en la ciudad fue de 1 millón. Es obvio que la conclusión es falsa, puesto que las cifras son extremas y el tamaño de la muestra es demasiado pequeño; por lo tanto la muestra no es representativa.

Hay muchas otras clases de datos inadecuados. Por ejemplo, algunos datos son respuestas inexactas de una encuesta, porque las preguntas usadas en la misma son vagas o engañosas, algunos datos son toscas estimaciones porque no hay disponibles datos exactos o es demasiado costosa su obtención, y algunos datos son irrelevantes en un problema dado, porque el estudio estadístico no está bien planeado. Al momento de recopilar los datos que serán procesados se es susceptible de cometer errores así como durante los cálculos de los mismos. No obstante, hay otros errores que no tienen nada que ver con la digitación y que no son tan fácilmente identificables. Algunos de éstos errores son:

- **Sesgo:** Es imposible ser completamente objetivo o no tener ideas preconcebidas antes de comenzar a estudiar un problema, y existen muchas maneras en que una perspectiva o estado mental pueda influir en la recopilación y en el análisis de la información. En estos casos se dice que hay un sesgo cuando el individuo da mayor peso a los datos que apoyan su opinión que a aquellos que la contradicen. Un caso extremo de sesgo sería la situación donde primero se toma una decisión y después se utiliza el análisis estadístico para justificar la decisión ya tomada.
- **Datos No Comparables:** el establecer comparaciones es una de las partes más importantes del análisis estadístico, pero es extremadamente importante que tales comparaciones se hagan entre datos que sean comparables.

- **Proyección descuidada de tendencias:** la proyección simplista de tendencias pasadas hacia el futuro es uno de los errores que más ha desacreditado el uso del análisis estadístico.
- **Muestreo Incorrecto:** en la mayoría de los estudios sucede que el volumen de información disponible es tan inmenso que se hace necesario estudiar muestras, para derivar conclusiones acerca de la población a que pertenece la muestra. Si la muestra se selecciona correctamente, tendrá básicamente las mismas propiedades que la población de la cual fue extraída; pero si el muestreo se realiza incorrectamente, entonces puede suceder que los resultados no signifiquen nada.

Sesgo significa que un usuario dé los datos perjudicialmente de más énfasis a los hechos, los cuales son empleados para mantener su predeterminada posición u opinión. Los estadísticos son frecuentemente degradados por lemas tales como: *Hay tres clases de mentiras: mentiras, mentiras reprobables y estadística, y Las cifras no mienten, pero los mentirosos piensan.* Hay dos clases de sesgos: conscientes e inconscientes. Ambos son comunes en el análisis estadístico. Hay numerosos ejemplos de sesgos conscientes. Un anunciante frecuentemente usa la estadística para probar que su producto es muy superior al producto de su competidor. Un político prefiere usar la estadística para sostener su punto de vista. Gerentes y líderes de trabajadores pueden simultáneamente situar sus respectivas cifras estadísticas sobre la misma tabla de trato para mostrar que sus rechazos o peticiones son justificadas. Es casi imposible que un sesgo inconsciente esté completamente ausente en un trabajo estadístico. En lo que respecta al ser humano, es difícil obtener una actitud completamente objetiva al abordar un problema, aun cuando un científico debería tener una mente abierta. Un estadístico debería estar enterado del hecho de que su interpretación de los resultados del análisis estadístico está influenciado por su propia experiencia, conocimiento y antecedentes con relación al problema dado.

3.2 Fundamentos

3.2.1 Pruebas de Hipótesis

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado.
- El interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.
- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

Una prueba de hipótesis está formada por cinco partes:

- La hipótesis nula, denotada por H_0 .
- La hipótesis alternativa, denotada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 3. Las dos hipótesis en competencia son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones:

- Rechazar H_0 y concluir que H_1 es verdadera.
- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 1. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19 usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 2. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas:

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 4. El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadístico de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**. Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**. Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente. Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 5. Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

3.2.2 Muestras grandes: una media poblacional

Definición 6. El valor de p (p -value) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.

Nota 10. Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.

Nota 11. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 7. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados:

p	H_0	Significativa
$p \leq 0.01$	rechazada	Result. altamente significativos y en contra de H_0
$p \leq 0.05$	rechazada	Result. Estadísticamente significativos y en contra de H_0
$p \leq 0.10$	rechazada	Result. posiblemente significativos con Tendencia estadística y en contra de H_0
$p > 0.10$	no rechazada	Result. estadísticamente no significativos y no rechazar H_0

Nota 12. Para determinar el valor de p , encontrar el área en la cola después del estadístico de prueba. Si la prueba es de una cola, este es el valor de p . Si es de dos colas, éste valor encontrado es la mitad del valor de p . Rechazar H_0 cuando el valor de $p < \alpha$.

Hay dos tipos de errores al realizar una prueba de hipótesis

	H_0 es Verdadera	H_0 es Falsa
Rechazar H_0	Error tipo I	✓
Aceptar H_0	✓	Error tipo II

Definición 8. La probabilidad de cometer el error tipo II se define por β donde

$$\begin{aligned}\beta &= P\{\text{error tipo II}\} = P\{\text{Aceptar equivocadamente } H_0\} \\ &= P\{\text{Aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\}\end{aligned}$$

Nota 13. Cuando H_0 es falsa y H_1 es verdadera, no siempre es posible especificar un valor exacto de μ , sino más bien un rango de posibles valores. En lugar de arriesgarse a tomar una decisión incorrecta, es mejor concluir que no hay evidencia suficiente para rechazar H_0 , es decir en lugar de aceptar H_0 , no rechazar H_0 .

La bondad de una prueba estadística se mide por el tamaño de α y β , ambas deben de ser pequeñas. Una manera muy efectiva de medir la potencia de la prueba es calculando el complemento del error tipo II:

$$\begin{aligned}1 - \beta &= P\{\text{Rechazar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{\text{Rechazar } H_0 \text{ cuando } H_1 \text{ es verdadera}\}\end{aligned}$$

Definición 9. La potencia de la prueba, $1 - \beta$, mide la capacidad de que la prueba funciona como se necesita.

Ejemplo 3. La producción diaria de una planta química local ha promediado 880 toneladas en los últimos años. A la gerente de control de calidad le gustaría saber si este promedio ha cambiado en meses recientes. Ella selecciona al azar 50 días de la base de datos computarizada y calcula el promedio y la desviación estándar de las $n = 50$ producciones como $\bar{x} = 871$ toneladas y $s = 21$ toneladas, respectivamente. Pruebe la hipótesis apropiada usando $\alpha = 0.05$.

La hipótesis nula apropiada es:

$$\begin{aligned} H_0 &: \mu = 880 \\ &\text{y la hipótesis alternativa } H_1 \text{ es} \\ H_1 &: \mu \neq 880 \end{aligned}$$

el estimador puntual para μ es \bar{x} , entonces el estadístico de prueba es

$$\begin{aligned} z &= \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \\ &= \frac{871 - 880}{21/\sqrt{50}} = -3.03 \end{aligned}$$

Para esta prueba de dos colas, hay que determinar los dos valores de $z_{\alpha/2}$, es decir, $z_{\alpha/2} = \pm 1.96$, como $z > z_{\alpha/2}$, z cae en la zona de rechazo, por lo tanto la gerente puede rechazar la hipótesis nula y concluir que el promedio efectivamente ha cambiado. La probabilidad de rechazar H_0 cuando esta es verdadera es de 0.05. Recordemos que el valor observado del estadístico de prueba es $z = -3.03$, la región de rechazo más pequeña que puede usarse y todavía seguir rechazando H_0 es $|z| > 3.03$, entonces $p = 2(0.012) = 0.0024$, que a su vez es menor que el nivel de significancia α asignado inicialmente, y además los resultados son **altamente significativos**. Finalmente determinemos la potencia de la prueba cuando μ en realidad es igual a 870 toneladas.

Recordar que la región de aceptación está entre -1.96 y 1.96 , para $\mu = 880$, equivalentemente

$$874.18 < \bar{x} < 885.82$$

β es la probabilidad de aceptar H_0 cuando $\mu = 870$, calculemos los valores de z correspondientes a 874.18 y 885.82. Entonces

$$\begin{aligned} z_1 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{874.18 - 870}{21/\sqrt{50}} = 1.41 \\ z_2 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{885.82 - 870}{21/\sqrt{50}} = 5.33 \end{aligned}$$

por lo tanto

$$\begin{aligned} \beta &= P\{\text{aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\} = P\{874.18 < \bar{x} < 885.82 \text{ cuando } \mu = 870\} \\ &= P\{1.41 < z < 5.33\} = P\{1.41 < z\} = 1 - 0.9207 = 0.0793 \end{aligned}$$

entonces, la potencia de la prueba es

$$1 - \beta = 1 - 0.0793 = 0.9207$$

que es la probabilidad de rechazar correctamente H_0 cuando H_0 es falsa.

Prueba de hipótesis para la diferencia entre dos medias poblacionales

El estadístico que resume la información muestral respecto a la diferencia en medias poblacionales ($\mu_1 - \mu_2$) es la diferencia de las medias muestrales ($\bar{x}_1 - \bar{x}_2$), por tanto al probar la diferencia entre las medias muestrales

se verifica que la diferencia real entre las medias poblacionales difiere de un valor especificado, $(\mu_1 - \mu_2) = D_0$, se puede usar el error estándar de $(\bar{x}_1 - \bar{x}_2)$, es decir

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \quad (3.1)$$

cuyo estimador está dado por

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} \quad (3.2)$$

El procedimiento para muestras grandes es:

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$, donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

- 2) **Hipótesis Alternativa**

Prueba de una Cola	Prueba de dos colas
$H_1 : (\mu_1 - \mu_2) > D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$
$H_1 : (\mu_1 - \mu_2) < D_0$	

- 3) Estadístico de prueba:

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (3.3)$$

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_0$	
$z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cuando $p < \alpha$	

Ejemplo 4. Para determinar si ser propietario de un automóvil afecta el rendimiento académico de un estudiante, se tomaron dos muestras aleatorias de 100 estudiantes varones. El promedio de calificaciones para los $n_1 = 100$ no propietarios de un auto tuvieron un promedio y varianza de $\bar{x}_1 = 2.7$ y $s_1^2 = 0.36$, respectivamente, mientras que para la segunda muestra con $n_2 = 100$ propietarios de un auto, se tiene $\bar{x}_2 = 2.54$ y $s_2^2 = 0.4$. Los datos presentan suficiente evidencia para indicar una diferencia en la media en el rendimiento académico entre propietarios y no propietarios de un automóvil? Hacer pruebas para $\alpha = 0.01, 0.05$ y $\alpha = 0.1$.

- Solución utilizando la técnica de regiones de rechazo: realizando las operaciones $z = 1.84$, determinar si excede los valores de $z_{\alpha/2}$.
- Solución utilizando el p-value: Calcular el valor de p , la probabilidad de que z sea mayor que $z = 1.84$ o menor que $z = -1.84$, se tiene que $p = 0.0658$.
- Si el intervalo de confianza que se construye contiene el valor del parámetro especificado por H_0 , entonces ese valor es uno de los posibles valores del parámetro y H_0 no debe ser rechazada.
- Si el valor hipotético se encuentra fuera de los límites de confianza, la hipótesis nula es rechazada al nivel de significancia α .

Prueba de Hipótesis para una Proporción Binomial

Para una muestra aleatoria de n intentos idénticos, de una población binomial, la proporción muestral \hat{p} tiene una distribución aproximadamente normal cuando n es grande, con media p y error estándar

$$SE = \sqrt{\frac{pq}{n}}. \quad (3.4)$$

La prueba de hipótesis de la forma

$$\begin{aligned} H_0 &: p = p_0 \\ H_1 &: p > p_0, \text{ o } p < p_0 \text{ o } p \neq p_0 \end{aligned}$$

El estadístico de prueba se construye con el mejor estimador de la proporción verdadera, \hat{p} , con el estadístico de prueba z , que se distribuye normal estándar.

El procedimiento es

- 1) Hipótesis nula: $H_0 : p = p_0$
- 2) Hipótesis alternativa

Prueba de una Cola	Prueba de dos colas
$H_1 : p > p_0$	$p \neq p_0$
$H_1 : p < p_0$	

- 3) Estadístico de prueba:

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{pq}{n}}}, \hat{p} = \frac{x}{n} \quad (3.5)$$

donde x es el número de éxitos en n intentos binomiales.

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_0$	
$z < -z_\alpha$ cuando $H_1 : p < p_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cuando $p < \alpha$	

Prueba de Hipótesis diferencia entre dos Proporciones Binomiales

Cuando se tienen dos muestras aleatorias independientes de dos poblaciones binomiales, el objetivo del experimento puede ser la diferencia $(p_1 - p_2)$ en las proporciones de individuos u objetos que poseen una característica específica en las dos poblaciones. En este caso se pueden utilizar los estimadores de las dos proporciones $(\hat{p}_1 - \hat{p}_2)$ con error estándar dado por

$$SE = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}, \quad (3.6)$$

considerando el estadístico z con un nivel de significancia $(1 - \alpha) 100\%$

La hipótesis nula a probarse es de la forma

$H_0: p_1 = p_2$ o equivalentemente $(p_1 - p_2) = 0$, contra una hipótesis alternativa H_1 de una o dos colas.

Para estimar el error estándar del estadístico z , se debe de utilizar el hecho de que suponiendo que H_0 es verdadera, las dos proporciones son iguales a algún valor común, p . Para obtener el mejor estimador de p es

$$p = \frac{\text{número total de éxitos}}{\text{Número total de pruebas}} = \frac{x_1 + x_2}{n_1 + n_2}. \quad (3.7)$$

1) **Hipótesis Nula:** $H_0 : (p_1 - p_2) = 0$

2) **Hipótesis Alternativa:** $H_1 :$

Prueba de una Cola	Prueba de dos colas
$H_1 : (p_1 - p_2) > 0$	$H_1 : (p_1 - p_2) \neq 0$
$H_1 : (p_1 - p_2) < 0$	

3) Estadístico de prueba:

$$z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{pq}{n_1} + \frac{pq}{n_2}}}, \quad (3.8)$$

donde $\hat{p}_1 = x_1/n_1$ y $\hat{p}_2 = x_2/n_2$, dado que el valor común para p_1 y p_2 es p , entonces $\hat{p} = \frac{x_1+x_2}{n_1+n_2}$ y por tanto el estadístico de prueba es

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}. \quad (3.9)$$

4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_\alpha$	
$z < -z_\alpha$ cuando $H_1 : p < p_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cuando $p < \alpha$	

3.2.3 Muestras Pequeñas

1) **Hipótesis Nula:** $H_0 : \mu = \mu_0$,

2) **Hipótesis Alternativa:** $H_1 :$

Prueba de una Cola	Prueba de dos colas
$H_1 : \mu > \mu_0$	$H_1 : \mu \neq \mu_0$
$H_1 : \mu < \mu_0$	

3) Estadístico de prueba:

$$t = \frac{\bar{x} - \mu_0}{\sqrt{\frac{s^2}{n}}}, \quad (3.10)$$

4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$t > t_\alpha$	
$t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$
cuando $p < \alpha$	

Diferencia entre dos medias poblacionales: MAI

Cuando los tamaños de muestra son pequeños, no se puede asegurar que las medias muestrales sean normales, pero si las poblaciones originales son normales, entonces la distribución muestral de la diferencia de las medias muestrales, $(\bar{x}_1 - \bar{x}_2)$, será normal con media $(\mu_1 - \mu_2)$ y error estándar

$$ES = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}. \quad (3.11)$$

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

- 2) **Hipótesis Alternativa**

Prueba de una Cola	Prueba de dos colas
$H_1 : (\mu_1 - \mu_2) > D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$
$H_1 : (\mu_1 - \mu_2) < D_0$	

- 3) Estadístico de prueba:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (3.12)$$

donde

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}. \quad (3.13)$$

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_0$	
$z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cuando $p < \alpha$	

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

Diferencia entre dos medias poblacionales: Diferencias Pareadas

- 1) **Hipótesis Nula:** $H_0 : \mu_d = 0$

- 2) **Hipótesis Alternativa:** $H_1 : \mu_d$

Prueba de una Cola	Prueba de dos colas
$H_1 : \mu_d > 0$	$H_1 : \mu_d \neq 0$
$H_1 : \mu_d < 0$	

- 3) Estadístico de prueba:

$$t = \frac{\bar{d}}{\sqrt{\frac{s_d^2}{n}}} \quad (3.14)$$

donde n es el número de diferencias pareadas, \bar{d} es la media de las diferencias muestrales, y s_d es la desviación estándar de las diferencias muestrales.

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

Inferencias con respecto a la Varianza Poblacional

- 1) **Hipótesis Nula:** $H_0 : \sigma^2 = \sigma_0^2$

- 2) **Hipótesis Alternativa:** H_1

Prueba de una Cola	Prueba de dos colas
$H_1 : \sigma^2 > \sigma_0^2$ $H_1 : \sigma^2 < \sigma_0^2$	$H_1 : \sigma^2 \neq \sigma_0^2$

- 3) Estadístico de prueba:

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}, \quad (3.15)$$

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$\chi^2 > \chi_\alpha^2$ $\chi^2 < \chi_{(1-\alpha)}^2$ cuando $H_1 : \chi^2 < \chi_0^2$ cuando $p < \alpha$	$\chi^2 > \chi_{\alpha/2}^2$ o $\chi^2 < \chi_{(1-\alpha/2)}^2$

Los valores críticos de χ^2 , están basados en $(n_1 +)$ grados de libertad.

Comparación de dos varianzas poblacionales

- 1) **Hipótesis Nula** $H_0 : (\sigma_1^2 - \sigma_2^2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

- 2) **Hipótesis Alternativa**

Prueba de una Cola	Prueba de dos colas
$H_1 : (\sigma_1^2 - \sigma_2^2) > D_0$ $H_1 : (\sigma_1^2 - \sigma_2^2) < D_0$	$H_1 : (\sigma_1^2 - \sigma_2^2) \neq D_0$

- 3) Estadístico de prueba:

$$F = \frac{s_1^2}{s_2^2} \quad (3.16)$$

donde s_1^2 es la varianza muestral más grande.

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$F > F_\alpha$ cuando $p < \alpha$	$F > F_{\alpha/2}$

3.2.4 Estimación por intervalos

Recordemos que S^2 es un estimador insesgado de σ^2 , entonces se tiene la siguiente definición

Definición 10. Sean $\hat{\theta}_1$ y $\hat{\theta}_2$ dos estimadores insesgados de θ , parámetro poblacional. Si $\sigma_{\hat{\theta}_1}^2 < \sigma_{\hat{\theta}_2}^2$, decimos que $\hat{\theta}_1$ es un estimador más eficaz de θ que $\hat{\theta}_2$.

Algunas observaciones que es preciso realizar

Nota 14. a) Para poblaciones normales, \bar{X} y \tilde{X} son estimadores insesgados de μ , pero con $\sigma_{\bar{X}}^2 < \sigma_{\tilde{X}}^2$.

b) Para las estimaciones por intervalos de θ , un intervalo de la forma $\hat{\theta}_L < \theta < \hat{\theta}_U$, $\hat{\theta}_L$ y $\hat{\theta}_U$ dependen del valor de $\hat{\theta}$.

c) Para $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$, si $n \rightarrow \infty$, entonces $\hat{\theta} \rightarrow \mu$.

Nota 15. Para $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$, si $n \rightarrow \infty$,

d) Para $\hat{\theta}$ se determinan $\hat{\theta}_L$ y $\hat{\theta}_U$ de modo tal que

$$P \left\{ \hat{\theta}_L < \hat{\theta} < \hat{\theta}_U \right\} = 1 - \alpha, \quad (3.17)$$

con $\alpha \in (0, 1)$. Es decir, $\theta \in (\hat{\theta}_L, \hat{\theta}_U)$ es un intervalo de confianza del $100(1 - \alpha)\%$.

e) De acuerdo con el TLC se espera que la distribución muestral de \bar{X} se distribuye aproximadamente normal con media $\mu_X = \mu$ y desviación estándar $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$.

Para $Z_{\alpha/2}$ se tiene $P \left\{ -Z_{\alpha/2} < Z < Z_{\alpha/2} \right\} = 1 - \alpha$, donde $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$. Entonces

$$P \left\{ -Z_{\alpha/2} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < Z_{\alpha/2} \right\} = 1 - \alpha, \quad (3.18)$$

es equivalente a

$$P \left\{ \bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right\} = 1 - \alpha. \quad (3.19)$$

f) Si \bar{X} es la media muestral de una muestra de tamaño n de una población con varianza conocida σ^2 , el intervalo de confianza de $100(1 - \alpha)\%$ para μ es

$$\mu \in \left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right). \quad (3.20)$$

g) Para muestras pequeñas de poblaciones no normales, no se puede esperar que el grado de confianza sea preciso.

h) Para $n \geq 30$, con distribución de forma no muy sesgada, se pueden tener buenos resultados.

Teorema 1. Si \bar{X} es un estimador de μ , podemos tener $100(1 - \alpha)\%$ de confianza en que el error no excederá a $z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$, error entre \bar{X} y μ .

Teorema 2. Si \bar{X} es un estimador de μ , podemos tener $100(1 - \alpha)\%$ de confianza en que el error no excederá una cantidad e cuando el tamaño de la muestra es

$$n = \left(\frac{z_{\alpha/2} \sigma}{e} \right)^2. \quad (3.21)$$

Nota 16. Para intervalos unilaterales

$$P \left\{ \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < Z_\alpha \right\} = 1 - \alpha \quad (3.22)$$

equivalentemente

$$P \left\{ \mu < \bar{X} + Z_\alpha \frac{\sigma}{\sqrt{n}} \right\} = 1 - \alpha. \quad (3.23)$$

Si \bar{X} es la media de una muestra aleatoria de tamaño n a partir de una población con varianza σ^2 , los límites de confianza unilaterales del $100(1 - \alpha)\%$ de confianza para μ están dados por

- a) Límite unilateral superior: $\bar{x} + z_\alpha \frac{\sigma}{\sqrt{n}}$
- b) Límite unilateral inferior: $\bar{x} - z_\alpha \frac{\sigma}{\sqrt{n}}$
- c) Para σ desconocida recordar que $T = \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_{n-1}$, donde s es la desviación estándar de la muestra. Entonces

$$P \left\{ -t_{\alpha/2} < T < t_{\alpha/2} \right\} = 1 - \alpha, \text{ equivalentemente} \quad (3.24)$$

$$P \left\{ \bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}} \right\} = 1 - \alpha. \quad (3.25)$$

- d) Un intervalo de confianza del $100(1 - \alpha)\%$ de confianza para μ , σ^2 desconocida y población normal es

$$\mu \in \left(\bar{x} - t_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}} \right) \quad (3.26)$$

,

donde $t_{\alpha/2}$ es una t -student con $\nu = n - 1$ grados de libertad.

- e) Los límites unilaterales para μ con σ desconocida son $\bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}}$ y $\bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}$.
- f) Cuando la población no es normal, σ desconocida y $n \geq 30$, σ se puede reemplazar por s para obtener el intervalo de confianza para muestras grandes:

$$\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}. \quad (3.27)$$

- g) El estimador de \bar{X} de μ , σ desconocida, la varianza de $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$, el error estándar de \bar{X} es σ/\sqrt{n} .

- h) Si σ es desconocida y la población es normal, $s \rightarrow \sigma$ y se incluye el error estándar s/\sqrt{n} , entonces

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}. \quad (3.28)$$

Intervalos de confianza sobre la varianza

Supongamos que X se distribuye normal (μ, σ^2) , desconocidas. Sea X_1, X_2, \dots, X_n muestra aleatoria de tamaño n , s^2 la varianza muestral.

Se sabe que $X^2 = \frac{(n-1)s^2}{\sigma^2}$ se distribuye χ_{n-1}^2 grados de libertad. Su intervalo de confianza es

$$\begin{aligned} P \left\{ \chi_{1-\frac{\alpha}{2}, n-1}^2 \leq \chi^2 \leq \chi_{\frac{\alpha}{2}, n-1}^2 \right\} &= 1 - \alpha \\ P \left\{ \chi_{1-\frac{\alpha}{2}, n-1}^2 \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi_{\frac{\alpha}{2}, n-1}^2 \right\} &= 1 - \alpha \\ P \left\{ \frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right\} &= 1 - \alpha, \end{aligned} \quad (3.29)$$

es decir

$$\sigma^2 \in \left[\frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2}, \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right], \quad (3.30)$$

los intervalos unilaterales son

$$\sigma^2 \in \left[\frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2}, \infty \right), \quad (3.31)$$

y

$$\sigma^2 \in \left(-\infty, \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right]. \quad (3.32)$$

Intervalos de confianza para proporciones

Supongamos que se tienen una muestra de tamaño n de una población grande pero finita, y supongamos que X , $X \leq n$, pertenecen a la clase de interés, entonces

$$\hat{p} = \frac{\bar{X}}{n}, \quad (3.33)$$

es el estimador puntual de la proporción de la población que pertenece a dicha clase. n y p son los parámetros de la distribución binomial, entonces

$$\hat{p} \sim N \left(p, \frac{p(1-p)}{n} \right) \quad (3.34)$$

aproximadamente si p es distinto de 0 y 1; o si n es suficientemente grande. Entonces

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim N(0, 1), \text{ aproximadamente.} \quad (3.35)$$

Entonces

$$\begin{aligned} 1 - \alpha &= P \left\{ -z_{\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \leq z_{\alpha/2} \right\} \\ &= P \left\{ \hat{p} - z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \right\} \end{aligned} \quad (3.36)$$

con $\sqrt{\frac{p(1-p)}{n}}$ error estándar del estimador puntual p . Una solución para determinar el intervalo de confianza del parámetro p (desconocido) es

$$1 - \alpha = P \left\{ \hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right\} \quad (3.37)$$

entonces los intervalos de confianza, tanto unilaterales como de dos colas son:

$$\text{a) } p \in \left(\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right),$$

$$\text{b) } p \in \left(-\infty, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right),$$

$$c) \quad p \in \left(\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \infty \right);$$

para minimizar el error estándar, se propone que el tamaño de la muestra sea

$$n = \left(\frac{z_{\alpha/2}}{E} \right)^2 p(1-p), \quad (3.38)$$

donde

$$E = |p - \hat{p}|.$$

Intervalos de confianza para dos muestras: Varianzas conocidas

Sean X_1 y X_2 variables aleatorias independientes. X_1 con media desconocida μ_1 y varianza conocida σ_1^2 ; y X_2 con media desconocida μ_2 y varianza conocida σ_2^2 . Se busca encontrar un intervalo de confianza de $100(1-\alpha)\%$ de la diferencia entre medias μ_1 y μ_2 . Sean $X_{11}, X_{12}, \dots, X_{1n_1}$ muestra aleatoria de n_1 observaciones de X_1 , y sean $X_{21}, X_{22}, \dots, X_{2n_2}$ muestra aleatoria de n_2 observaciones de X_2 .

Sean \bar{X}_1 y \bar{X}_2 , medias muestrales, entonces el estadístico

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1), \quad (3.39)$$

si X_1 y X_2 son normales o aproximadamente normales si se aplican las condiciones del Teorema de Límite Central respectivamente. Entonces se tiene

$$1 - \alpha = P\{-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}\} = P\left\{-Z_{\alpha/2} \leq \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \leq Z_{\alpha/2}\right\} \quad (3.40)$$

$$= P\left\{(\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \leq \mu_1 - \mu_2 \leq (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right\}. \quad (3.41)$$

Entonces los intervalos de confianza unilaterales y de dos colas al $(1-\alpha)\%$ de confianza son

a)

$$\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right], \quad (3.42)$$

b)

$$\mu_1 - \mu_2 \in \left[-\infty, (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right], \quad (3.43)$$

c)

$$\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \infty \right]. \quad (3.44)$$

Nota 17. Si σ_1 y σ_2 son conocidas, o por lo menos se conoce una aproximación, y los tamaños de las muestras n_1 y n_2 son iguales, $n_1 = n_2 = n$, se puede determinar el tamaño de la muestra para que el error al estimar $\mu_1 - \mu_2$ usando $\bar{X}_1 - \bar{X}_2$ sea menor que E (valor del error deseado) al $(1-\alpha)\%$ de confianza. El tamaño n de la muestra requerido para cada muestra es

$$n = \left(\frac{Z_{\alpha/2}}{E} \right)^2 (\sigma_1^2 + \sigma_2^2). \quad (3.45)$$

Intervalos de confianza para dos muestras: Varianzas desconocidas e iguales

- a) Si $n_1, n_2 \geq 30$ se pueden utilizar los intervalos de la distribución normal para varianzas conocidas
- b) Si n_1, n_2 son muestras pequeñas, supongase que las poblaciones para X_1 y X_2 son normales con varianzas desconocidas y con base en el intervalo de confianza para distribuciones t -student

Supongamos que X_1 es una variable aleatoria con media μ_1 y varianza σ_1^2 , X_2 es una variable aleatoria con media μ_2 y varianza σ_2^2 . Todos los parámetros son desconocidos. Sin embargo supóngase que es razonable considerar que $\sigma_1^2 = \sigma_2^2 = \sigma^2$.

Nuevamente sean X_1 y X_2 variables aleatorias independientes. X_1 con media desconocida μ_1 y varianza muestral S_1^2 ; y X_2 con media desconocida μ_2 y varianza muestral S_2^2 . Dado que S_1^2 y S_2^2 son estimadores de σ_1^2 , se propone el estimador S de σ^2 como

$$S_p^2 = \frac{(n_1 - 1) S_1^2 + (n_2 - 1) S_2^2}{n_1 + n_2 - 2}, \quad (3.46)$$

entonces, el estadístico para $\mu_1 - \mu_2$ es

$$t_\nu = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}, \quad (3.47)$$

donde t_ν es una t de student con $\nu = n_1 + n_2 - 2$ grados de libertad.

Por lo tanto

$$\begin{aligned} 1 - \alpha &= P \{ -t_{\alpha/2, \nu} \leq t \leq t_{\alpha/2, \nu} \} \\ &= P \left\{ (\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \leq t \leq (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right\}, \end{aligned} \quad (3.48)$$

luego, los intervalos de confianza del $(1 - \alpha) \%$ para $\mu_1 - \mu_2$ son

a)

$$\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right]. \quad (3.49)$$

b)

$$\mu_1 - \mu_2 \in \left[-\infty, (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right]. \quad (3.50)$$

c)

$$\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \infty \right]. \quad (3.51)$$

Intervalos de confianza para dos muestras: Varianzas desconocidas diferentes

Si no se tiene certeza de que $\sigma_1^2 = \sigma_2^2$, se propone el estadístico

$$t^* = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}, \quad (3.52)$$

que se distribuye t -student con ν grados de libertad, donde

$$\nu = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\frac{S_1^2/n_1}{n_1+1} + \frac{S_2^2/n_2}{n_2+1}} - 2. \quad (3.53)$$

Entonces el intervalo de confianza de aproximadamente el $100(1 - \alpha)\%$ para $\mu_1 - \mu_2$ con $\sigma_1^2 \neq \sigma_2^2$ es

$$\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}, (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \right]. \quad (3.54)$$

Intervalos de confianza para razón de Varianzas

Supongamos que se toman dos muestras aleatorias independientes de las dos poblaciones de interés. Sean X_1 y X_2 variables normales independientes con medias desconocidas μ_1 y μ_2 y varianzas desconocidas σ_1^2 y σ_2^2 respectivamente. Se busca un intervalo de confianza de $100(1 - \alpha)\%$ para σ_1^2/σ_2^2 . Supongamos n_1 y n_2 muestras aleatorias de X_1 y X_2 y sean S_1^2 y S_2^2 varianzas muestrales. Se sabe que

$$F = \frac{S_2^2/\sigma_2^2}{S_1^2/\sigma_1^2}, \quad (3.55)$$

se distribuye F con $n_2 - 1$ y $n_1 - 1$ grados de libertad.

Por lo tanto

$$\begin{aligned} P \left\{ F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq F \leq F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right\} &= 1 - \alpha, \\ P \left\{ F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq \frac{S_2^2/\sigma_2^2}{S_1^2/\sigma_1^2} \leq F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right\} &= 1 - \alpha, \end{aligned} \quad (3.56)$$

luego entonces

$$P \left\{ \frac{S_1^2}{S_2^2} F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{S_1^2}{S_2^2} F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right\} = 1 - \alpha. \quad (3.57)$$

en consecuencia

$$\frac{\sigma_1^2}{\sigma_2^2} \in \left[\frac{S_1^2}{S_2^2} F_{1-\frac{\alpha}{2}, n_2-1, n_1-1}, \frac{S_1^2}{S_2^2} F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right], \quad (3.58)$$

donde

$$F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} = \frac{1}{F_{\frac{\alpha}{2}, n_2-1, n_1-1}}. \quad (3.59)$$

Intervalos de confianza para diferencia de proporciones

Sean dos proporciones de interés p_1 y p_2 . Se busca un intervalo para $p_1 - p_2$ al $100(1 - \alpha)\%$. Sean dos muestras independientes de tamaño n_1 y n_2 de poblaciones infinitas de modo que X_1 y X_2 variables aleatorias binomiales independientes con parámetros (n_1, p_1) y (n_2, p_2) . X_1 y X_2 son el número de observaciones que pertenecen a la clase de interés correspondientes. Entonces $\hat{p}_1 = \frac{X_1}{n_1}$ y $\hat{p}_2 = \frac{X_2}{n_2}$ son estimadores de p_1 y p_2 respectivamente. Supongamos que se cumple la aproximación normal a la binomial, entonces

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}} \sim N(0, 1) \text{ aproximadamente} \quad (3.60)$$

por tanto

$$(\hat{p}_1 - \hat{p}_2) - Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \leq p_1 - p_2 \leq (\hat{p}_1 - \hat{p}_2) + Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \quad (3.61)$$

3.3 Análisis de Regresión Lineal (RL)

En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**. Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n . La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente. A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

3.3.1 Regresión Lineal Simple (RLS)

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria. El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x,$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (3.62)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*. Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación realiza por el **métodos de mínimos cuadrados**.

Entonces la ecuación (??) se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (3.63)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2. \quad (3.64)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\frac{\partial L}{\partial \beta_0} = 0, \quad \frac{\partial L}{\partial \beta_1} = 0.$$

Evalutando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0, \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0, \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i, \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i. \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \quad (3.65)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2}, \quad (3.66)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x. \quad (3.67)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2, \quad (3.68)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right); \quad (3.69)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}. \quad (3.70)$$

3.3.2 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo. Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$\begin{aligned}
E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x})\right) \\
&= \frac{1}{S_{xx}} \left[\beta_0 E\left(\sum_{k=1}^n (x_k - \bar{x})\right) + E\left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x})\right) + E\left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x})\right) \right] \\
&= \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1.
\end{aligned}$$

Por lo tanto

$$E(\hat{\beta}_1) = \beta_1, \quad (3.71)$$

Es decir, $\hat{\beta}_1$ es un estimador insesgado. Ahora calculemos la varianza:

$$\begin{aligned}
V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) = \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) \\
&= \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}},
\end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}. \quad (3.72)$$

Entonces tenemos la siguiente proposición:

Proposición 5.

$$E(\hat{\beta}_0) = \beta_0, \quad (3.73)$$

$$V(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \quad (3.74)$$

$$Cov(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\sigma^2 \bar{x}}{S_{xx}}. \quad (3.75)$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$$e_k = y_k - \hat{y}_k, \text{ se le denomina } \mathbf{residuo}.$$

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*:

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (3.76)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$SC_E = \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \quad (3.77)$$

$$E(SC_E) = (n-2)\sigma^2, \text{ por lo tanto,} \quad (3.78)$$

$$\hat{\sigma}^2 = \frac{SC_E}{n-2} = \mathbf{MC_E} \text{ es un estimador insesgado de } \sigma^2. \quad (3.79)$$

3.3.3 Prueba de Hipótesis en RLS

Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza. Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$. Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_1, \sigma^2)$. De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??). Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (3.80)$$

que se distribuye t con $n - 2$ grados de libertad bajo $H_0 : \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (3.81)$$

Para β_0 se puede proceder de manera análoga para

$$H_0 : \beta_0 = \beta_{0,0},$$

$$H_1 : \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}, \quad (3.82)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (3.83)$$

- No rechazar $H_0 : \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0 : \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0 : \beta_1 = 0$ puede realizarse de la siguiente manera:

$$\begin{aligned} S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 = \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\ &= \sum_{k=1}^n \left[(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2 \right] \\ &= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \end{aligned}$$

$$\begin{aligned}
& \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\
& = \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
& + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
& + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0.
\end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (3.84)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (3.85)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (3.86)$$

Por lo tanto la ecuación (??) se puede reescribir como:

$$S_{yy} = SC_R + SC_E, \quad (3.87)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$:

$$\begin{aligned}
S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}), \\
S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R.
\end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 6.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx}, \quad (3.88)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E},$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$. El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n - 2$	MC_E	
Total	S_{yy}	$n - 1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir,

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}}. \quad (3.89)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}.$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

3.3.4 Estimación de Intervalos en RLS

Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros. El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión. Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n - 2$ grados de libertad. Por tanto un intervalo de confianza de $100(1 - \alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}}. \quad (3.90)$$

De igual manera, para β_0 un intervalo de confianza al $100(1 - \alpha)\%$ es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (3.91)$$

3.3.5 Predicción

Supongamos que se tiene un valor x_0 de interés, entonces la estimación puntual de este nuevo valor

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0. \quad (3.92)$$

Esta nueva observación es independiente de las utilizadas para obtener el modelo de regresión, por tanto, el intervalo en torno a la recta de regresión es inapropiado, puesto que se basa únicamente en los datos empleados para ajustar el modelo de regresión. El intervalo de confianza en torno a la recta de regresión se refiere a la respuesta media verdadera $x = x_0$, no a observaciones futuras. Sea y_0 la observación futura en $x = x_0$, y sea \hat{y}_0 dada en la ecuación anterior, el estimador de y_0 . Si se define la variable aleatoria

$$w = y_0 - \hat{y}_0,$$

esta se distribuye normalmente con media cero y varianza

$$V(w) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right] \quad (3.93)$$

dado que y_0 es independiente de \hat{y}_0 , por lo tanto el intervalo de predicción al nivel α para futuras observaciones x_0 es

$$\hat{y}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]} \leq y_0 \leq \hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]}. \quad (3.94)$$

3.3.6 Prueba de falta de ajuste

Es común encontrar que el modelo ajustado no satisface totalmente el modelo necesario para los datos, en este caso es preciso saber qué tan bueno es el modelo propuesto. Para esto se propone la siguiente prueba de hipótesis:

H_0 : El modelo propuesto se ajusta adecuadamente a los datos.

H_1 : El modelo NO se ajusta a los datos.

La prueba implica dividir la suma de cuadrados del error o del residuo en las siguientes dos componentes:

$$SC_E = SC_{EP} + SC_{FDA} \quad (3.95)$$

donde SC_{EP} es la suma de cuadrados atribuibles al error puro, y SC_{FDA} es la suma de cuadrados atribuible a la falta de ajuste del modelo.

3.3.7 Coeficiente de Determinación

La cantidad

$$R^2 = \frac{SC_R}{S_{yy}} = 1 - \frac{SC_E}{S_{yy}}, \quad (3.96)$$

se denomina coeficiente de determinación y se utiliza para saber si el modelo de regresión es suficiente o no. Se puede demostrar que $0 \leq R^2 \leq 1$, una manera de interpretar este valor es que si $R^2 = k$, entonces el modelo de regresión explica el $k * 100\%$ de la variabilidad en los datos. R^2 . Este coeficiente tiene las siguientes propiedades

- No mide la magnitud de la pendiente de la recta de regresión.
- Un valor grande de R^2 no implica una pendiente empinada.
- No mide la suficiencia del modelo.
- Valores grandes de R^2 no implican necesariamente que el modelo de regresión proporcionará predicciones precisas para futuras observaciones.

Chapter 4

Regresión Logística

4.1 Introducción

La regresión logística es una técnica de modelado estadístico ampliamente utilizada en análisis de datos cuando el objetivo es predecir la probabilidad de un resultado binario, es decir, cuando la variable dependiente o respuesta tiene dos posibles categorías, como "éxito/fallo" o "sí/no". Esta técnica se emplea en una variedad de disciplinas, como la biomedicina, ciencias sociales, marketing y más, para resolver problemas donde la variable respuesta es discreta o categórica.

A diferencia de la regresión lineal, que asume una relación lineal entre las variables independientes y la variable dependiente y que produce valores en un rango continuo, la regresión logística está diseñada para manejar situaciones donde la respuesta es categórica. En su forma más común, la regresión logística binaria, el modelo predice la probabilidad de que un evento ocurra en función de una o más variables independientes. Este tipo de regresión toma la forma de un modelo no lineal, debido a la naturaleza discreta de la variable dependiente.

La regresión lineal busca modelar la relación entre una variable dependiente continua Y y una o más variables independientes X_1, X_2, \dots, X_n mediante una ecuación de la forma:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon, \quad (4.1)$$

donde $\beta_0, \beta_1, \dots, \beta_n$ son los coeficientes del modelo y ϵ es el término de error. La regresión logística, en cambio, modela la probabilidad de que un evento ocurra (por ejemplo, éxito vs. fracaso) utilizando la función logística. La variable dependiente Y es binaria, tomando valores de 0 o 1. La ecuación de la regresión logística es:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n, \quad (4.2)$$

donde p es la probabilidad de que $Y = 1$. La función logística es:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}. \quad (4.3)$$

La regresión logística se utiliza en una variedad de campos para problemas de clasificación binaria, tales como:

- **Medicina:** Predicción de la presencia o ausencia de una enfermedad.
- **Marketing:** Determinación de la probabilidad de que un cliente compre un producto.
- **Finanzas:** Evaluación del riesgo de crédito, es decir, si un cliente va a incumplir o no con un préstamo.
- **Seguridad:** Detección de fraudes o intrusiones.

4.1.1 Implementación Básica en R

Para implementar una regresión logística en R, primero es necesario instalar y cargar los paquetes necesarios. Aquí se muestra un ejemplo básico de implementación:

- Descargue e instale R desde <https://cran.r-project.org/>.
- Descargue e instale RStudio desde <https://rstudio.com/products/rstudio/download/>.

```
# Instalación del paquete necesario
install.packages("stats")

# Carga del paquete
library(stats)

# Ejemplo de conjunto de datos
data <- data.frame(
  outcome = c(1, 0, 1, 0, 1, 1, 0, 1, 0, 0),
  predictor = c(2.3, 1.9, 3.1, 2.8, 3.6, 2.4, 2.1, 3.3, 2.2, 1.7)
)

# Ajuste del modelo de regresión logística
model <- glm(outcome ~ predictor, data = data, family = binomial)

# Resumen del modelo
summary(model)
```

4.2 Conceptos Básicos

La regresión logística es una técnica de modelado estadístico utilizada para predecir la probabilidad de un evento binario (es decir, un evento que tiene dos posibles resultados) en función de una o más variables independientes. Un modelo de regresión logística describe cómo una variable dependiente binaria Y (que puede tomar los valores 0 o 1) está relacionada con una o más variables independientes X_1, X_2, \dots, X_n . A diferencia de la regresión lineal, que predice un valor continuo, la regresión logística predice una probabilidad que puede ser interpretada como la probabilidad de que $Y = 1$ dado un conjunto de valores para X_1, X_2, \dots, X_n .

4.2.1 Regresión Lineal

La regresión lineal es utilizada para predecir el valor de una variable dependiente continua en función de una o más variables independientes. El modelo de regresión lineal tiene la forma:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon \quad (4.4)$$

donde:

- Y es la variable dependiente.
- β_0 es la intersección con el eje Y o término constante.
- $\beta_1, \beta_2, \dots, \beta_n$ son los coeficientes que representan la relación entre las variables independientes y la variable dependiente.
- X_1, X_2, \dots, X_n son las variables independientes.

- e) ϵ es el término de error, que representa la desviación de los datos observados de los valores predichos por el modelo.

El objetivo de la regresión lineal es encontrar los valores de los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ que minimicen la suma de los cuadrados de las diferencias entre los valores observados y los valores predichos. Este método se conoce como mínimos cuadrados ordinarios (OLS, por sus siglas en inglés). La función de costo a minimizar es:

$$J(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (4.5)$$

donde:

- a) y_i es el valor observado de la variable dependiente para la i -ésima observación.
b) \hat{y}_i es el valor predicho por el modelo para la i -ésima observación, dado por:

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_n x_{in} \quad (4.6)$$

Para encontrar los valores óptimos de los coeficientes, se toman las derivadas parciales de la función de costo con respecto a cada coeficiente y se igualan a cero:

$$\frac{\partial J}{\partial \beta_j} = 0 \quad \text{para } j = 0, 1, \dots, n \quad (4.7)$$

Resolviendo este sistema de ecuaciones, se obtienen los valores de los coeficientes que minimizan la función de costo.

4.2.2 Regresión Logística

La deducción de la fórmula de la regresión logística comienza con la necesidad de modelar la probabilidad de un evento binario. Queremos encontrar una función que relacione las variables independientes con la probabilidad de que la variable dependiente tome el valor 1. La probabilidad de que el evento ocurra, $P(Y = 1)$, se denota como p . La probabilidad de que el evento no ocurra, $P(Y = 0)$, es $1 - p$. Los **odds** (chances) de que ocurra el evento se definen como:

$$\text{odds} = \frac{p}{1 - p} \quad (4.8)$$

Los odds indican cuántas veces más probable es que ocurra el evento frente a que no ocurra. Para simplificar el modelado de los *odds*, se aplica el logaritmo natural, obteniendo la función **logit**:

$$\text{logit}(p) = \log\left(\frac{p}{1 - p}\right) \quad (4.9)$$

La transformación logit es útil porque convierte el rango de la probabilidad $(0, 1)$ al rango de números reales $(-\infty, \infty)$. La idea clave de la regresión logística es modelar la transformación logit de la probabilidad como una combinación lineal de las variables independientes:

$$\text{logit}(p) = \log\left(\frac{p}{1 - p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad (4.10)$$

Aquí, β_0 es el término constante y $\beta_1, \beta_2, \dots, \beta_n$ son los coeficientes asociados con las variables independientes X_1, X_2, \dots, X_n . Para expresar p en función de una combinación lineal de las variables independientes, invertimos la transformación logit. Partimos de la ecuación ??, aplicando la función exponencial en ambos lados:

$$\frac{p}{1 - p} = e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n} \quad (4.11)$$

Despejando p :

$$p = \frac{e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n}}{1 + e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n}} \quad (4.12)$$

La expresión final que obtenemos es conocida como la **función logística**:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} \quad (4.13)$$

Esta función describe cómo las variables independientes se relacionan con la probabilidad de que el evento de interés ocurra. Los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ se estiman a partir de los datos utilizando el método de máxima verosimilitud.

4.3 Método de Máxima Verosimilitud

Para estimar los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ en la regresión logística, utilizamos el método de máxima verosimilitud. La idea es encontrar los valores de los coeficientes que maximicen la probabilidad de observar los datos dados. Esta probabilidad se expresa mediante la función de verosimilitud L . La función de verosimilitud $L(\beta_0, \beta_1, \dots, \beta_n)$ para un conjunto de n observaciones se define como el producto de las probabilidades de las observaciones dadas las variables independientes:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (4.14)$$

donde:

- p_i es la probabilidad predicha de que $Y_i = 1$,
- y_i es el valor observado de la variable dependiente para la i -ésima observación.

Trabajar directamente con esta función de verosimilitud puede ser complicado debido al producto de muchas probabilidades, especialmente si n es grande. Para simplificar los cálculos, se utiliza el logaritmo de la función de verosimilitud, conocido como la función de log-verosimilitud. El uso del logaritmo simplifica significativamente la diferenciación y maximización de la función. La función de log-verosimilitud se define como:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (4.15)$$

Aquí, \log representa el logaritmo natural. Esta transformación es válida porque el logaritmo es una función monótona creciente, lo que significa que maximizar la log-verosimilitud es equivalente a maximizar la verosimilitud original. En la regresión logística, la probabilidad p_i está dada por la función logística:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \quad (4.16)$$

Sustituyendo esta expresión en la función de log-verosimilitud, obtenemos:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n \left[y_i \log \left(\frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \right) + \right. \\ &\quad \left. (1 - y_i) \log \left(1 - \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \right) \right] \end{aligned}$$

Simplificando esta expresión, notamos que:

$$\log\left(\frac{1}{1+e^{-z}}\right) = -\log(1+e^{-z})$$

y

$$\log\left(1 - \frac{1}{1+e^{-z}}\right) = \log\left(\frac{e^{-z}}{1+e^{-z}}\right) = -z - \log(1+e^{-z})$$

Aplicando estas identidades, la función de log-verosimilitud se convierte en:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n \left[y_i (-\log(1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})})) + \right. \\ &\quad \left. (1 - y_i) \left(-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}) \right) \right] \end{aligned}$$

Simplificando aún más, obtenemos:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n [y_i(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) \\ &\quad - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}})] \end{aligned}$$

Para simplificar aún más la notación, podemos utilizar notación matricial. Definimos la matriz \mathbf{X} de tamaño $n \times (k+1)$ y el vector de coeficientes $\boldsymbol{\beta}$ de tamaño $(k+1) \times 1$ como sigue:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \quad (4.17)$$

Entonces, la expresión para la función de log-verosimilitud es:

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n \left[y_i (\mathbf{X}_i \boldsymbol{\beta}) - \log(1 + e^{\mathbf{X}_i \boldsymbol{\beta}}) \right] \quad (4.18)$$

donde \mathbf{X}_i es la i -ésima fila de la matriz \mathbf{X} . Esta notación matricial simplifica la implementación y la derivación de los estimadores de los coeficientes en la regresión logística. Utilizando métodos numéricos, como el algoritmo de Newton-Raphson, se pueden encontrar los coeficientes que maximizan la función de log-verosimilitud. Para maximizar la función de log-verosimilitud, derivamos esta función con respecto a cada uno de los coeficientes β_j y encontramos los puntos críticos. La derivada parcial de la función de log-verosimilitud con respecto a β_j es:

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = \sum_{i=1}^n \left[y_i X_{ij} - \frac{X_{ij} e^{\mathbf{X}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{X}_i \boldsymbol{\beta}}} \right] \quad (4.19)$$

Simplificando, esta derivada se puede expresar como:

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = \sum_{i=1}^n X_{ij} (y_i - p_i), \quad \text{donde } p_i = \frac{1}{1 + e^{-\mathbf{X}_i \boldsymbol{\beta}}} \quad (4.20)$$

Para encontrar los coeficientes que maximizan la log-verosimilitud, resolvemos el sistema de ecuaciones

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = 0 \text{ para todos los } j = 0, 1, \dots, k.$$

Este sistema de ecuaciones no tiene una solución analítica cerrada, por lo que se resuelve numéricamente utilizando métodos iterativos como el algoritmo de Newton-Raphson.

4.4 Método de Newton-Raphson

El método de Newton-Raphson es un algoritmo iterativo que se utiliza para encontrar las raíces de una función. En el contexto de la regresión logística, se utiliza para maximizar la función de log-verosimilitud encontrando los valores de los coeficientes $\beta_0, \beta_1, \dots, \beta_n$. Este método se basa en una aproximación de segundo orden de la función objetivo. Dado un valor inicial de los coeficientes $\boldsymbol{\beta}^{(0)}$, se actualiza iterativamente el valor de los coeficientes utilizando la fórmula:

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - [\mathbf{H}(\boldsymbol{\beta}^{(t)})]^{-1} \nabla \log L(\boldsymbol{\beta}^{(t)}) \quad (4.21)$$

donde:

- $\boldsymbol{\beta}^{(t)}$ es el vector de coeficientes en la t -ésima iteración.
- $\nabla \log L(\boldsymbol{\beta}^{(t)})$ es el gradiente de la función de log-verosimilitud con respecto a los coeficientes $\boldsymbol{\beta}$:

$$\nabla \log L(\boldsymbol{\beta}) = \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \quad (4.22)$$

donde \mathbf{y} es el vector de valores observados y \mathbf{p} es el vector de probabilidades.

- $\mathbf{H}(\boldsymbol{\beta}^{(t)})$ es la matriz Hessiana (matriz de segundas derivadas) evaluada en $\boldsymbol{\beta}^{(t)}$:

$$\mathbf{H}(\boldsymbol{\beta}) = -\mathbf{X}^T \mathbf{W} \mathbf{X} \quad (4.23)$$

donde \mathbf{W} es una matriz diagonal de pesos con elementos $w_i = p_i(1 - p_i)$.

En resumen:

Algoritmo 1. *El algoritmo Newton-Raphson para la regresión logística se puede resumir en los siguientes pasos:*

1. *Inicializar el vector de coeficientes $\boldsymbol{\beta}^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).*
2. *Calcular el gradiente $\nabla \log L(\boldsymbol{\beta}^{(t)})$ y la matriz Hessiana $\mathbf{H}(\boldsymbol{\beta}^{(t)})$ en la iteración t .*
3. *Actualizar los coeficientes utilizando la fórmula:*

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - [\mathbf{H}(\boldsymbol{\beta}^{(t)})]^{-1} \nabla \log L(\boldsymbol{\beta}^{(t)}) \quad (4.24)$$

4. *Repetir los pasos 2 y 3 hasta que la diferencia entre $\boldsymbol{\beta}^{(t+1)}$ y $\boldsymbol{\beta}^{(t)}$ sea menor que un umbral predefinido (criterio de convergencia).*

En resumen, el método de Newton-Raphson permite encontrar los coeficientes que maximizan la función de log-verosimilitud de manera eficiente.

4.5 Especificando

En específico para un conjunto de n observaciones, la función de verosimilitud L se define como el producto de las probabilidades individuales de observar cada dato:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (4.25)$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación y p_i es la probabilidad predicha de que $Y_i = 1$. Aquí, p_i es dado por la función logística:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \quad (4.26)$$

Tomando el logaritmo:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (4.27)$$

Sustituyendo p_i :

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}) \right] \quad (4.28)$$

Dado que el objetivo es encontrar los valores de $\beta_0, \beta_1, \dots, \beta_n$ que maximicen la función de log-verosimilitud. Para β_j , la derivada parcial de la función de log-verosimilitud es:

$$\frac{\partial \log L}{\partial \beta_j} = \sum_{i=1}^n \left[y_i X_{ij} - \frac{X_{ij} e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}}{1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}} \right] \quad (4.29)$$

Esto se simplifica a (comparar con la ecuación ??):

$$\frac{\partial \log L}{\partial \beta_j} = \sum_{i=1}^n X_{ij} (y_i - p_i) \quad (4.30)$$

Para maximizar la log-verosimilitud, resolvemos el sistema de ecuaciones $\frac{\partial \log L}{\partial \beta_j} = 0$ para todos los j de 0 a n , mismo que se resuelve numéricamente utilizando métodos el algoritmo de Newton-Raphson. El método de Newton-Raphson se basa en una aproximación de segundo orden de la función objetivo. Dado un valor inicial de los coeficientes $\beta^{(0)}$, se iterativamente actualiza el valor de los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - \left[\mathbf{H}(\beta^{(k)}) \right]^{-1} \mathbf{g}(\beta^{(k)}) \quad (4.31)$$

donde:

- $\beta^{(k)}$ es el vector de coeficientes en la k -ésima iteración.
- $\mathbf{g}(\beta^{(k)})$ es el gradiente (vector de primeras derivadas) evaluado en $\beta^{(k)}$:

$$\mathbf{g}(\beta) = \frac{\partial \log L}{\partial \beta} = \sum_{i=1}^n \mathbf{X}_i (y_i - p_i) \quad (4.32)$$

donde \mathbf{X}_i es el vector de valores de las variables independientes para la i -ésima observación (comparar con ecuación ??).

- $\mathbf{H}(\beta^{(k)})$ es la matriz Hessiana (matriz de segundas derivadas) evaluada en $\beta^{(k)}$:

$$\mathbf{H}(\beta) = \frac{\partial^2 \log L}{\partial \beta \partial \beta^T} = - \sum_{i=1}^n p_i(1-p_i) \mathbf{X}_i \mathbf{X}_i^T, \quad (4.33)$$

comparar con ecuación ??

Algoritmo 2. Los pasos del algoritmo Newton-Raphson para la regresión logística son:

1. Inicializar el vector de coeficientes $\beta^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).
2. Calcular el gradiente $\mathbf{g}(\beta^{(k)})$ y la matriz Hessiana $\mathbf{H}(\beta^{(k)})$ en la iteración k .
3. Actualizar los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - [\mathbf{H}(\beta^{(k)})]^{-1} \mathbf{g}(\beta^{(k)}) \quad (4.34)$$

4. Repetir los pasos 2 y 3 hasta que la diferencia entre $\beta^{(k+1)}$ y $\beta^{(k)}$ sea menor que un umbral predefinido (criterio de convergencia).

Como se puede observar la diferencia entre el Algoritmo ?? y el Algoritmo ?? son mínimas

Notas finales

En el contexto de la regresión logística, los vectores X_1, X_2, \dots, X_n representan las variables independientes. Cada X_j es un vector columna que contiene los valores de la variable independiente j para cada una de las n observaciones. Es decir,

$$X_j = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{bmatrix} \quad (4.35)$$

Para simplificar la notación y los cálculos, a menudo combinamos todos los vectores de variables independientes en una única matriz de diseño \mathbf{X} de tamaño $n \times (k+1)$, donde n es el número de observaciones y $k+1$ es el número de variables independientes más el término de intercepto. La primera columna de \mathbf{X} corresponde a un vector de unos para el término de intercepto, y las demás columnas corresponden a los valores de las variables independientes:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix} \quad (4.36)$$

revisar la ecuación ??. De esta forma, el modelo logit puede ser escrito de manera compacta utilizando la notación matricial:

$$\text{logit}(p) = \log \left(\frac{p}{1-p} \right) = \mathbf{X}\beta \quad (4.37)$$

donde β es el vector de coeficientes:

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \quad (4.38)$$

Así, la probabilidad p se puede expresar como:

$$p = \frac{1}{1 + e^{-\mathbf{X}\boldsymbol{\beta}}} \quad (4.39)$$

Comparar la ecuación anterior con la ecuación ???. Esta notación matricial simplifica la implementación y la derivación de los estimadores de los coeficientes en la regresión logística. Para estimar los coeficientes $\boldsymbol{\beta}$ en la regresión logística, se utiliza el método de máxima verosimilitud. La función de verosimilitud $L(\boldsymbol{\beta})$ se define como el producto de las probabilidades de las observaciones dadas las variables independientes, recordemos la ecuación ??:

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (4.40)$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación, y p_i es la probabilidad predicha de que $Y_i = 1$. La función de log-verosimilitud, que es más fácil de maximizar, se obtiene tomando el logaritmo natural de la función de verosimilitud (??):

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (4.41)$$

Sustituyendo $p_i = \frac{1}{1+e^{-\mathbf{X}_i\boldsymbol{\beta}}}$, donde \mathbf{X}_i es la i -ésima fila de la matriz de diseño \mathbf{X} , obtenemos:

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n \left[y_i (\mathbf{X}_i\boldsymbol{\beta}) - \log(1 + e^{\mathbf{X}_i\boldsymbol{\beta}}) \right] \quad (4.42)$$

Para encontrar los valores de $\boldsymbol{\beta}$ que maximizan la función de log-verosimilitud, se utiliza un algoritmo iterativo como el método de Newton-Raphson. Este método requiere calcular el gradiente y la matriz Hessiana de la función de log-verosimilitud.

El gradiente de la función de log-verosimilitud con respecto a $\boldsymbol{\beta}$ es (?? y ??):

$$\nabla \log L(\boldsymbol{\beta}) = \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \quad (4.43)$$

donde \mathbf{y} es el vector de valores observados y \mathbf{p} es el vector de probabilidades predichas.

La matriz Hessiana de la función de log-verosimilitud es (?? y ??):

$$\mathbf{H}(\boldsymbol{\beta}) = -\mathbf{X}^T \mathbf{W} \mathbf{X} \quad (4.44)$$

donde \mathbf{W} es una matriz diagonal de pesos con elementos $w_i = p_i(1 - p_i)$.

El método de Newton-Raphson actualiza los coeficientes $\boldsymbol{\beta}$ de la siguiente manera:

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - [\mathbf{H}(\boldsymbol{\beta}^{(t)})]^{-1} \nabla \log L(\boldsymbol{\beta}^{(t)}) \quad (4.45)$$

Iterando este proceso hasta que la diferencia entre $\boldsymbol{\beta}^{(t+1)}$ y $\boldsymbol{\beta}^{(t)}$ sea menor que un umbral predefinido (??, ??, ?? y ??), se obtienen los estimadores de máxima verosimilitud para los coeficientes de la regresión logística.

Chapter 5

Elementos de Probabilidad

5.1 Introducción

Los fundamentos de probabilidad y estadística son esenciales para comprender y aplicar técnicas de análisis de datos y modelado estadístico, incluyendo la regresión lineal y logística. Este capítulo proporciona una revisión de los conceptos clave en probabilidad y estadística que son relevantes para estos métodos.

5.2 Probabilidad

La probabilidad es una medida de la incertidumbre o el grado de creencia en la ocurrencia de un evento. Los conceptos fundamentales incluyen:

5.2.1 Espacio Muestral y Eventos

El espacio muestral, denotado como S , es el conjunto de todos los posibles resultados de un experimento aleatorio. Un evento es un subconjunto del espacio muestral. Por ejemplo, si lanzamos un dado, el espacio muestral es:

$$S = \{1, 2, 3, 4, 5, 6\}$$

Un evento podría ser obtener un número par:

$$E = \{2, 4, 6\}$$

5.2.2 Definiciones de Probabilidad

Existen varias definiciones de probabilidad, incluyendo la probabilidad clásica, la probabilidad frecuentista y la probabilidad bayesiana.

Probabilidad Clásica

La probabilidad clásica se define como el número de resultados favorables dividido por el número total de resultados posibles:

$$P(E) = \frac{|E|}{|S|}$$

donde $|E|$ es el número de elementos en el evento E y $|S|$ es el número de elementos en el espacio muestral S .

Probabilidad Frecuentista

La probabilidad frecuentista se basa en la frecuencia relativa de ocurrencia de un evento en un gran número de repeticiones del experimento:

$$P(E) = \lim_{n \rightarrow \infty} \frac{n_E}{n}$$

donde n_E es el número de veces que ocurre el evento E y n es el número total de repeticiones del experimento.

Probabilidad Bayesiana

La probabilidad bayesiana se interpreta como un grado de creencia actualizado a medida que se dispone de nueva información. Se basa en el teorema de Bayes:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

donde $P(A|B)$ es la probabilidad de A dado B , $P(B|A)$ es la probabilidad de B dado A , $P(A)$ y $P(B)$ son las probabilidades de A y B respectivamente.

5.3 Estadística Bayesiana

La estadística bayesiana proporciona un enfoque coherente para el análisis de datos basado en el teorema de Bayes. Los conceptos fundamentales incluyen:

5.3.1 Prior y Posterior

Distribución Prior

La distribución prior (apriori) representa nuestra creencia sobre los parámetros antes de observar los datos. Es una distribución de probabilidad que refleja nuestra incertidumbre inicial sobre los parámetros. Por ejemplo, si creemos que un parámetro θ sigue una distribución normal con media μ_0 y varianza σ_0^2 , nuestra prior sería:

$$P(\theta) = \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(\theta-\mu_0)^2}{2\sigma_0^2}}$$

Verosimilitud

La verosimilitud (likelihood) es la probabilidad de observar los datos dados los parámetros. Es una función de los parámetros θ dada una muestra de datos X :

$$L(\theta; X) = P(X|\theta)$$

donde X son los datos observados y θ son los parámetros del modelo.

Distribución Posterior

La distribución posterior (a posteriori) combina la información de la prior y la verosimilitud utilizando el teorema de Bayes. Representa nuestra creencia sobre los parámetros después de observar los datos:

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

donde $P(\theta|X)$ es la distribución posterior, $P(X|\theta)$ es la verosimilitud, $P(\theta)$ es la prior y $P(X)$ es la probabilidad marginal de los datos.

La probabilidad marginal de los datos $P(X)$ se puede calcular como:

$$P(X) = \int_{\Theta} P(X|\theta)P(\theta)d\theta$$

donde Θ es el espacio de todos los posibles valores del parámetro θ .

5.4 Distribuciones de Probabilidad

Las distribuciones de probabilidad describen cómo se distribuyen los valores de una variable aleatoria. Existen distribuciones de probabilidad discretas y continuas.

5.4.1 Distribuciones Discretas

Una variable aleatoria discreta toma un número finito o contable de valores. Algunas distribuciones discretas comunes incluyen:

Distribución Binomial

La distribución binomial describe el número de éxitos en una serie de ensayos de Bernoulli independientes y con la misma probabilidad de éxito. La función de probabilidad es:

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

donde X es el número de éxitos, n es el número de ensayos, p es la probabilidad de éxito en cada ensayo, y $\binom{n}{k}$ es el coeficiente binomial.

La función generadora de momentos (MGF) para la distribución binomial es:

$$M_X(t) = (1 - p + pe^t)^n$$

El valor esperado y la varianza de una variable aleatoria binomial son:

$$\begin{aligned} E(X) &= np \\ \text{Var}(X) &= np(1 - p) \end{aligned}$$

Distribución de Poisson

La distribución de Poisson describe el número de eventos que ocurren en un intervalo de tiempo fijo o en un área fija. La función de probabilidad es:

$$P(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

donde X es el número de eventos, λ es la tasa media de eventos por intervalo, y k es el número de eventos observados.

La función generadora de momentos (MGF) para la distribución de Poisson es:

$$M_X(t) = e^{\lambda(e^t - 1)}$$

El valor esperado y la varianza de una variable aleatoria de Poisson son:

$$\begin{aligned} E(X) &= \lambda \\ \text{Var}(X) &= \lambda \end{aligned}$$

5.4.2 Distribuciones Continuas

Una variable aleatoria continua toma un número infinito de valores en un intervalo continuo. Algunas distribuciones continuas comunes incluyen:

Distribución Normal

La distribución normal, también conocida como distribución gaussiana, es una de las distribuciones más importantes en estadística. La función de densidad de probabilidad es:

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

donde x es un valor de la variable aleatoria, μ es la media, y σ es la desviación estándar.

La función generadora de momentos (MGF) para la distribución normal es:

$$M_X(t) = e^{\mu t + \frac{1}{2}\sigma^2 t^2}$$

El valor esperado y la varianza de una variable aleatoria normal son:

$$\begin{aligned} E(X) &= \mu \\ \text{Var}(X) &= \sigma^2 \end{aligned}$$

Distribución Exponencial

La distribución exponencial describe el tiempo entre eventos en un proceso de Poisson. La función de densidad de probabilidad es:

$$f(x) = \lambda e^{-\lambda x}$$

donde x es el tiempo entre eventos y λ es la tasa media de eventos.

La función generadora de momentos (MGF) para la distribución exponencial es:

$$M_X(t) = \frac{\lambda}{\lambda - t}, \quad \text{para } t < \lambda$$

El valor esperado y la varianza de una variable aleatoria exponencial son:

$$\begin{aligned} E(X) &= \frac{1}{\lambda} \\ \text{Var}(X) &= \frac{1}{\lambda^2} \end{aligned}$$

5.5 Estadística Descriptiva

La estadística descriptiva resume y describe las características de un conjunto de datos. Incluye medidas de tendencia central, medidas de dispersión y medidas de forma.

5.5.1 Medidas de Tendencia Central

Las medidas de tendencia central incluyen la media, la mediana y la moda.

Media

La media aritmética es la suma de los valores dividida por el número de valores:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

donde x_i son los valores de la muestra y n es el tamaño de la muestra.

Mediana

La mediana es el valor medio cuando los datos están ordenados. Si el número de valores es impar, la mediana es el valor central. Si es par, es el promedio de los dos valores centrales.

Moda

La moda es el valor que ocurre con mayor frecuencia en un conjunto de datos.

5.5.2 Medidas de Dispersión

Las medidas de dispersión incluyen el rango, la varianza y la desviación estándar.

Rango

El rango es la diferencia entre el valor máximo y el valor mínimo de los datos:

$$Rango = x_{\max} - x_{\min}$$

Varianza

La varianza es la media de los cuadrados de las diferencias entre los valores y la media:

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Desviación Estándar

La desviación estándar es la raíz cuadrada de la varianza:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

5.6 Inferencia Estadística

La inferencia estadística es el proceso de sacar conclusiones sobre una población a partir de una muestra. Incluye la estimación de parámetros y la prueba de hipótesis.

5.6.1 Estimación de Parámetros

La estimación de parámetros implica el uso de datos muestrales para estimar los parámetros de una población.

Estimador Puntual

Un estimador puntual proporciona un único valor como estimación de un parámetro de la población. Por ejemplo, la media muestral \bar{x} es un estimador puntual de la media poblacional μ . Otros ejemplos de estimadores puntuales son:

- **Mediana muestral** (\tilde{x}): Estimador de la mediana poblacional.
- **Varianza muestral** (s^2): Estimador de la varianza poblacional σ^2 , definido como:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

- **Desviación estándar muestral** (s): Estimador de la desviación estándar poblacional σ , definido como:

$$s = \sqrt{s^2}$$

Propiedades de los Estimadores Puntuales

Los estimadores puntuales deben cumplir ciertas propiedades deseables, como:

- **Insesgadez**: Un estimador es insesgado si su valor esperado es igual al valor del parámetro que estima.

$$E(\hat{\theta}) = \theta$$

- **Consistencia**: Un estimador es consistente si converge en probabilidad al valor del parámetro a medida que el tamaño de la muestra tiende a infinito.
- **Eficiencia**: Un estimador es eficiente si tiene la varianza más baja entre todos los estimadores insesgados.

Estimador por Intervalo

Un estimador por intervalo proporciona un rango de valores dentro del cual se espera que se encuentre el parámetro poblacional con un cierto nivel de confianza. Por ejemplo, un intervalo de confianza para la media es:

$$\left(\bar{x} - z \frac{\sigma}{\sqrt{n}}, \bar{x} + z \frac{\sigma}{\sqrt{n}} \right)$$

donde z es el valor crítico correspondiente al nivel de confianza deseado, σ es la desviación estándar poblacional y n es el tamaño de la muestra.

5.6.2 Prueba de Hipótesis

La prueba de hipótesis es un procedimiento para decidir si una afirmación sobre un parámetro poblacional es consistente con los datos muestrales.

Hipótesis Nula y Alternativa

La hipótesis nula (H_0) es la afirmación que se somete a prueba, y la hipótesis alternativa (H_a) es la afirmación que se acepta si se rechaza la hipótesis nula.

Nivel de Significancia

El nivel de significancia (α) es la probabilidad de rechazar la hipótesis nula cuando es verdadera. Un valor comúnmente utilizado es $\alpha = 0.05$.

Estadístico de Prueba

El estadístico de prueba es una medida calculada a partir de los datos muestrales que se utiliza para decidir si se rechaza la hipótesis nula. Por ejemplo, en una prueba t para la media:

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

donde \bar{x} es la media muestral, μ_0 es la media poblacional bajo la hipótesis nula, s es la desviación estándar muestral y n es el tamaño de la muestra.

P-valor

El p-valor es la probabilidad de obtener un valor del estadístico de prueba al menos tan extremo como el observado, bajo la suposición de que la hipótesis nula es verdadera. Si el p-valor es menor que el nivel de significancia α , se rechaza la hipótesis nula. El p-valor se interpreta de la siguiente manera:

- **P-valor bajo ($p \leq 0.05$):** Evidencia suficiente para rechazar la hipótesis nula.
- **P-valor alto ($p > 0.05$):** No hay suficiente evidencia para rechazar la hipótesis nula.

Tipos de Errores

En la prueba de hipótesis, se pueden cometer dos tipos de errores:

- **Error Tipo I (α):** Rechazar la hipótesis nula cuando es verdadera.
- **Error Tipo II (β):** No rechazar la hipótesis nula cuando es falsa.

Tabla de Errores en la Prueba de Hipótesis

A continuación se presenta una tabla que muestra los posibles resultados en una prueba de hipótesis, incluyendo los falsos positivos (error tipo I) y los falsos negativos (error tipo II):

	Hipótesis Nula Verdadera	Hipótesis Nula Falsa
Rechazar H_0	Error Tipo I (α)	Aceptar H_a
No Rechazar H_0	Aceptar H_0	Error Tipo II (β)

Table 5.1: Resultados de la Prueba de Hipótesis

Chapter 6

Matemáticas Detrás de la Regresión Logística

6.1 Introducción

La regresión logística es una técnica de modelado estadístico utilizada para predecir la probabilidad de un evento binario en función de una o más variables independientes. Este capítulo profundiza en las matemáticas subyacentes a la regresión logística, incluyendo la función logística, la función de verosimilitud, y los métodos para estimar los coeficientes del modelo.

6.2 Función Logística

La función logística es la base de la regresión logística. Esta función transforma una combinación lineal de variables independientes en una probabilidad.

6.2.1 Definición

La función logística se define como:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$

donde p es la probabilidad de que el evento ocurra, $\beta_0, \beta_1, \dots, \beta_n$ son los coeficientes del modelo, y X_1, X_2, \dots, X_n son las variables independientes.

6.2.2 Propiedades

La función logística tiene varias propiedades importantes:

- **Rango:** La función logística siempre produce un valor entre 0 y 1, lo que la hace adecuada para modelar probabilidades.
- **Monotonía:** La función es monótona creciente, lo que significa que a medida que la combinación lineal de variables independientes aumenta, la probabilidad también aumenta.
- **Simetría:** La función logística es simétrica en torno a $p = 0.5$.

6.3 Función de Verosimilitud

La función de verosimilitud se utiliza para estimar los coeficientes del modelo de regresión logística. Esta función mide la probabilidad de observar los datos dados los coeficientes del modelo.

6.3.1 Definición

Para un conjunto de n observaciones, la función de verosimilitud L se define como el producto de las probabilidades individuales de observar cada dato:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación y p_i es la probabilidad predicha de que $Y_i = 1$.

6.3.2 Función de Log-Verosimilitud

Para simplificar los cálculos, trabajamos con el logaritmo de la función de verosimilitud, conocido como la función de log-verosimilitud. Tomar el logaritmo convierte el producto en una suma:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

Sustituyendo p_i :

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}) \right]$$

6.4 Estimación de Coeficientes

Los coeficientes del modelo de regresión logística se estiman maximizando la función de log-verosimilitud. Este proceso generalmente se realiza mediante métodos iterativos como el algoritmo de Newton-Raphson.

6.4.1 Gradiente y Hessiana

Para maximizar la función de log-verosimilitud, necesitamos calcular su gradiente y su matriz Hessiana.

Gradiente

El gradiente de la función de log-verosimilitud con respecto a los coeficientes β es:

$$\mathbf{g}(\beta) = \frac{\partial \log L}{\partial \beta} = \sum_{i=1}^n \mathbf{X}_i (y_i - p_i)$$

donde \mathbf{X}_i es el vector de valores de las variables independientes para la i -ésima observación.

Hessiana

La matriz Hessiana de la función de log-verosimilitud con respecto a los coeficientes β es:

$$\mathbf{H}(\beta) = \frac{\partial^2 \log L}{\partial \beta \partial \beta^T} = - \sum_{i=1}^n p_i (1 - p_i) \mathbf{X}_i \mathbf{X}_i^T$$

6.4.2 Algoritmo Newton-Raphson

El algoritmo Newton-Raphson se utiliza para encontrar los valores de los coeficientes que maximizan la función de log-verosimilitud. El algoritmo se puede resumir en los siguientes pasos:

1. Inicializar el vector de coeficientes $\beta^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).
2. Calcular el gradiente $\mathbf{g}(\beta^{(k)})$ y la matriz Hessiana $\mathbf{H}(\beta^{(k)})$ en la iteración k .
3. Actualizar los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - \left[\mathbf{H}(\beta^{(k)}) \right]^{-1} \mathbf{g}(\beta^{(k)})$$

4. Repetir los pasos 2 y 3 hasta que la diferencia entre $\beta^{(k+1)}$ y $\beta^{(k)}$ sea menor que un umbral predefinido (criterio de convergencia).

6.5 Validación del Modelo

Una vez que se han estimado los coeficientes del modelo de regresión logística, es importante validar el modelo para asegurarse de que proporciona predicciones precisas.

6.5.1 Curva ROC y AUC

La curva ROC (Receiver Operating Characteristic) es una herramienta gráfica utilizada para evaluar el rendimiento de un modelo de clasificación binaria. El área bajo la curva (AUC) mide la capacidad del modelo para distinguir entre las clases.

6.5.2 Matriz de Confusión

La matriz de confusión es una tabla que resume el rendimiento de un modelo de clasificación al comparar las predicciones del modelo con los valores reales. Los términos en la matriz de confusión incluyen verdaderos positivos, falsos positivos, verdaderos negativos y falsos negativos.

Chapter 7

Preparación de Datos y Selección de Variables

7.1 Introducción

La preparación de datos y la selección de variables son pasos cruciales en el proceso de modelado estadístico. Un modelo bien preparado y con las variables adecuadas puede mejorar significativamente la precisión y la interpretabilidad del modelo. Este capítulo proporciona una revisión detallada de las técnicas de limpieza de datos, tratamiento de datos faltantes, codificación de variables categóricas y selección de variables.

7.2 Importancia de la Preparación de Datos

La calidad de los datos es fundamental para el éxito de cualquier análisis estadístico. Los datos sin limpiar pueden llevar a modelos inexactos y conclusiones erróneas. La preparación de datos incluye varias etapas:

- Limpieza de datos
- Tratamiento de datos faltantes
- Codificación de variables categóricas
- Selección y transformación de variables

7.3 Limpieza de Datos

La limpieza de datos es el proceso de detectar y corregir (o eliminar) los datos incorrectos, incompletos o irrelevantes. Este proceso incluye:

- Eliminación de duplicados
- Corrección de errores tipográficos
- Consistencia de formato
- Tratamiento de valores extremos (outliers)

7.4 Tratamiento de Datos Faltantes

Los datos faltantes son un problema común en los conjuntos de datos y pueden afectar la calidad de los modelos. Hay varias estrategias para manejar los datos faltantes:

- **Eliminación de Datos Faltantes:** Se eliminan las filas o columnas con datos faltantes.
- **Imputación:** Se reemplazan los valores faltantes con estimaciones, como la media, la mediana o la moda.
- **Modelos Predictivos:** Se utilizan modelos predictivos para estimar los valores faltantes.

7.4.1 Imputación de la Media

Una técnica común es reemplazar los valores faltantes con la media de la variable. Esto se puede hacer de la siguiente manera:

$$x_i = \begin{cases} x_i & \text{si } x_i \text{ no es faltante} \\ \bar{x} & \text{si } x_i \text{ es faltante} \end{cases}$$

donde \bar{x} es la media de la variable.

7.5 Codificación de Variables Categóricas

Las variables categóricas deben ser convertidas a un formato numérico antes de ser usadas en un modelo de regresión logística. Hay varias técnicas para codificar variables categóricas:

7.5.1 Codificación One-Hot

La codificación one-hot crea una columna binaria para cada categoría. Por ejemplo, si tenemos una variable categórica con tres categorías (A, B, C), se crean tres columnas:

$$\begin{aligned} A &= [1, 0, 0] \\ B &= [0, 1, 0] \\ C &= [0, 0, 1] \end{aligned}$$

7.5.2 Codificación Ordinal

La codificación ordinal asigna un valor entero único a cada categoría, preservando el orden natural de las categorías. Por ejemplo:

$$\begin{aligned} \text{Bajo} &= 1 \\ \text{Medio} &= 2 \\ \text{Alto} &= 3 \end{aligned}$$

7.6 Selección de Variables

La selección de variables es el proceso de elegir las variables más relevantes para el modelo. Existen varias técnicas para la selección de variables:

7.6.1 Métodos de Filtrado

Los métodos de filtrado seleccionan variables basadas en criterios estadísticos, como la correlación o la chi-cuadrado. Algunas técnicas comunes incluyen:

- **Análisis de Correlación:** Se seleccionan variables con alta correlación con la variable dependiente y baja correlación entre ellas.
- **Pruebas de Chi-cuadrado:** Se utilizan para variables categóricas para determinar la asociación entre la variable independiente y la variable dependiente.

7.6.2 Métodos de Wrapper

Los métodos de wrapper evalúan múltiples combinaciones de variables y seleccionan la combinación que optimiza el rendimiento del modelo. Ejemplos incluyen:

- **Selección hacia Adelante:** Comienza con un modelo vacío y agrega variables una por una, seleccionando la variable que mejora más el modelo en cada paso.
- **Selección hacia Atrás:** Comienza con todas las variables y elimina una por una, removiendo la variable que tiene el menor impacto en el modelo en cada paso.
- **Selección Paso a Paso:** Combina la selección hacia adelante y hacia atrás, agregando y eliminando variables según sea necesario.

7.6.3 Métodos Basados en Modelos

Los métodos basados en modelos utilizan técnicas de regularización como Lasso y Ridge para seleccionar variables. Estas técnicas añaden un término de penalización a la función de costo para evitar el sobreajuste.

Regresión Lasso

La regresión Lasso (Least Absolute Shrinkage and Selection Operator) añade una penalización L_1 a la función de costo:

$$J(\beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

donde λ es el parámetro de regularización que controla la cantidad de penalización.

Regresión Ridge

La regresión Ridge añade una penalización L_2 a la función de costo:

$$J(\beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

donde λ es el parámetro de regularización.

7.7 Implementación en R

7.7.1 Limpieza de Datos

Para ilustrar la limpieza de datos en R, considere el siguiente conjunto de datos:

```
data <- data.frame(
  var1 = c(1, 2, 3, NA, 5),
  var2 = c("A", "B", "A", "B", "A"),
  var3 = c(10, 15, 10, 20, 25)
)

# Eliminación de filas con datos faltantes
data_clean <- na.omit(data)

# Imputación de la media
data$var1[is.na(data$var1)] <- mean(data$var1, na.rm = TRUE)
```

7.7.2 Codificación de Variables Categóricas

Para codificar variables categóricas, utilice la función ‘model.matrix’:

```
data <- data.frame(
  var1 = c(1, 2, 3, 4, 5),
  var2 = c("A", "B", "A", "B", "A")
)

# Codificación one-hot
data_onehot <- model.matrix(~ var2 - 1, data = data)
```

7.7.3 Selección de Variables

Para la selección de variables, utilice el paquete ‘caret’:

```
library(caret)

# Dividir los datos en conjuntos de entrenamiento y prueba
set.seed(123)
trainIndex <- createDataPartition(data$var1, p = .8,
                                   list = FALSE,
                                   times = 1)

dataTrain <- data[trainIndex,]
dataTest <- data[-trainIndex,]

# Modelo de regresión logística
model <- train(var1 ~ ., data = dataTrain, method = "glm", family = "binomial")

# Selección de variables
model <- stepAIC(model, direction = "both")
summary(model)
```

Chapter 8

Evaluación del Modelo y Validación Cruzada

8.1 Introducción

Evaluar la calidad y el rendimiento de un modelo de regresión logística es crucial para asegurar que las predicciones sean precisas y útiles. Este capítulo se centra en las técnicas y métricas utilizadas para evaluar modelos de clasificación binaria, así como en la validación cruzada, una técnica para evaluar la generalización del modelo.

8.2 Métricas de Evaluación del Modelo

Las métricas de evaluación permiten cuantificar la precisión y el rendimiento de un modelo. Algunas de las métricas más comunes incluyen:

8.2.1 Curva ROC y AUC

La curva ROC (Receiver Operating Characteristic) es una representación gráfica de la sensibilidad (verdaderos positivos) frente a 1 - especificidad (falsos positivos). El área bajo la curva (AUC) mide la capacidad del modelo para distinguir entre las clases.

$$\begin{aligned}\text{Sensibilidad} &= \frac{\text{TP}}{\text{TP} + \text{FN}} \\ \text{Especificidad} &= \frac{\text{TN}}{\text{TN} + \text{FP}}\end{aligned}$$

8.2.2 Matriz de Confusión

La matriz de confusión es una tabla que muestra el rendimiento del modelo comparando las predicciones con los valores reales. Los términos incluyen:

- **Verdaderos Positivos (TP):** Predicciones correctas de la clase positiva.
- **Falsos Positivos (FP):** Predicciones incorrectas de la clase positiva.
- **Verdaderos Negativos (TN):** Predicciones correctas de la clase negativa.
- **Falsos Negativos (FN):** Predicciones incorrectas de la clase negativa.

	Predicción Positiva	Predicción Negativa
Real Positiva	TP	FN
Real Negativa	FP	TN

Table 8.1: Matriz de Confusión

8.2.3 Precisión, Recall y F1-Score

$$\begin{aligned}\text{Precisión} &= \frac{\text{TP}}{\text{TP} + \text{FP}} \\ \text{Recall} &= \frac{\text{TP}}{\text{TP} + \text{FN}} \\ \text{F1-Score} &= 2 \cdot \frac{\text{Precisión} \cdot \text{Recall}}{\text{Precisión} + \text{Recall}}\end{aligned}$$

8.2.4 Log-Loss

La pérdida logarítmica (Log-Loss) mide la precisión de las probabilidades predichas. La fórmula es:

$$\text{Log-Loss} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

donde y_i son los valores reales y p_i son las probabilidades predichas.

8.3 Validación Cruzada

La validación cruzada es una técnica para evaluar la capacidad de generalización de un modelo. Existen varios tipos de validación cruzada:

8.3.1 K-Fold Cross-Validation

En K-Fold Cross-Validation, los datos se dividen en K subconjuntos. El modelo se entrena K veces, cada vez utilizando K-1 subconjuntos para el entrenamiento y el subconjunto restante para la validación.

$$\text{Error Medio} = \frac{1}{K} \sum_{k=1}^K \text{Error}_k$$

8.3.2 Leave-One-Out Cross-Validation (LOOCV)

En LOOCV, cada observación se usa una vez como conjunto de validación y las restantes como conjunto de entrenamiento. Este método es computacionalmente costoso pero útil para conjuntos de datos pequeños.

8.4 Ajuste y Sobreajuste del Modelo

El ajuste adecuado del modelo es crucial para evitar el sobreajuste (overfitting) y el subajuste (underfitting).

8.4.1 Sobreajuste

El sobreajuste ocurre cuando un modelo se ajusta demasiado bien a los datos de entrenamiento, capturando ruido y patrones irrelevantes. Los síntomas incluyen una alta precisión en el entrenamiento y baja precisión en la validación.

8.4.2 Subajuste

El subajuste ocurre cuando un modelo no captura los patrones subyacentes de los datos. Los síntomas incluyen baja precisión tanto en el entrenamiento como en la validación.

8.4.3 Regularización

La regularización es una técnica para prevenir el sobreajuste añadiendo un término de penalización a la función de costo. Las técnicas comunes incluyen:

- Regresión Lasso (L1)
- Regresión Ridge (L2)

8.5 Implementación en R

8.5.1 Evaluación del Modelo

```
# Cargar el paquete necesario
library(caret)

# Dividir los datos en conjuntos de entrenamiento y prueba
set.seed(123)
trainIndex <- createDataPartition(data$var1, p = .8,
                                   list = FALSE,
                                   times = 1)

dataTrain <- data[trainIndex,]
dataTest <- data[-trainIndex,]

# Entrenar el modelo de regresión logística
model <- train(var1 ~ ., data = dataTrain, method = "glm", family = "binomial")

# Predicciones en el conjunto de prueba
predictions <- predict(model, dataTest)

# Matriz de confusión
confusionMatrix(predictions, dataTest$var1)
```

8.5.2 Validación Cruzada

```
# K-Fold Cross-Validation
control <- trainControl(method = "cv", number = 10)
model_cv <- train(var1 ~ ., data = dataTrain, method = "glm",
                  family = "binomial", trControl = control)
```

```
# Evaluación del modelo con validación cruzada  
print(model_cv)
```

Chapter 9

Diagnóstico del Modelo y Ajuste de Parámetros

9.1 Introducción

El diagnóstico del modelo y el ajuste de parámetros son pasos esenciales para mejorar la precisión y la robustez de los modelos de regresión logística. Este capítulo se enfoca en las técnicas para diagnosticar problemas en los modelos y en métodos para ajustar los parámetros de manera óptima.

9.2 Diagnóstico del Modelo

El diagnóstico del modelo implica evaluar el rendimiento del modelo y detectar posibles problemas, como el sobreajuste, la multicolinealidad y la influencia de puntos de datos individuales.

9.2.1 Residuos

Los residuos son las diferencias entre los valores observados y los valores predichos por el modelo. El análisis de residuos puede revelar patrones que indican problemas con el modelo.

$$\text{Residuo}_i = y_i - \hat{y}_i$$

Residuos Estudiantizados

Los residuos estudiantizados se ajustan por la variabilidad del residuo y se utilizan para detectar outliers.

$$r_i = \frac{\text{Residuo}_i}{\hat{\sigma}\sqrt{1 - h_i}}$$

donde h_i es el leverage del punto de datos.

9.2.2 Influencia

La influencia mide el impacto de un punto de datos en los coeficientes del modelo. Los puntos con alta influencia pueden distorsionar el modelo.

Distancia de Cook

La distancia de Cook es una medida de la influencia de un punto de datos en los coeficientes del modelo.

$$D_i = \frac{r_i^2}{p} \cdot \frac{h_i}{1 - h_i}$$

donde p es el número de parámetros en el modelo.

9.2.3 Multicolinealidad

La multicolinealidad ocurre cuando dos o más variables independientes están altamente correlacionadas. Esto puede inflar las varianzas de los coeficientes y hacer que el modelo sea inestable.

Factor de Inflación de la Varianza (VIF)

El VIF mide cuánto se inflan las varianzas de los coeficientes debido a la multicolinealidad.

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

donde R_j^2 es el coeficiente de determinación de la regresión de la variable j contra todas las demás variables.

9.3 Ajuste de Parámetros

El ajuste de parámetros implica seleccionar los valores óptimos para los hiperparámetros del modelo. Esto puede mejorar el rendimiento y prevenir el sobreajuste.

9.3.1 Grid Search

El grid search es un método exhaustivo para ajustar los parámetros. Se define una rejilla de posibles valores de parámetros y se evalúa el rendimiento del modelo para cada combinación.

9.3.2 Random Search

El random search selecciona aleatoriamente combinaciones de valores de parámetros dentro de un rango especificado. Es menos exhaustivo que el grid search, pero puede ser más eficiente.

9.3.3 Bayesian Optimization

La optimización bayesiana utiliza modelos probabilísticos para seleccionar iterativamente los valores de parámetros más prometedores.

9.4 Implementación en R

9.4.1 Diagnóstico del Modelo

```
# Cargar el paquete necesario
library(car)
```

```

# Residuos estudentizados
dataTrain$resid <- rstudent(model)
hist(dataTrain$resid, breaks = 20, main = "Residuos Estudentizados")

# Distancia de Cook
dataTrain$cook <- cooks.distance(model)
plot(dataTrain$cook, type = "h", main = "Distancia de Cook")

# Factor de Inflaci\'on de la Varianza
vif_values <- vif(model)
print(vif_values)

```

9.4.2 Ajuste de Parámetros

```

# Grid Search con caret
control <- trainControl(method = "cv", number = 10)
tune_grid <- expand.grid(.alpha = c(0, 0.5, 1), .lambda = seq(0.01, 0.1, by = 0.01))

model_tune <- train(var1 ~ ., data = dataTrain, method = "glmnet",
                    trControl = control, tuneGrid = tune_grid)

print(model_tune)

```

Chapter 10

Interpretación de los Resultados

10.1 Introducción

Interpretar correctamente los resultados de un modelo de regresión logística es esencial para tomar decisiones informadas. Este capítulo se centra en la interpretación de los coeficientes del modelo, las odds ratios, los intervalos de confianza y la significancia estadística.

10.2 Coeficientes de Regresión Logística

Los coeficientes de regresión logística representan la relación entre las variables independientes y la variable dependiente en términos de log-odds.

10.2.1 Interpretación de los Coeficientes

Cada coeficiente β_j en el modelo de regresión logística se interpreta como el cambio en el log-odds de la variable dependiente por unidad de cambio en la variable independiente X_j .

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

10.2.2 Signo de los Coeficientes

- **Coeficiente Positivo:** Un coeficiente positivo indica que un aumento en la variable independiente está asociado con un aumento en el log-odds de la variable dependiente.
- **Coeficiente Negativo:** Un coeficiente negativo indica que un aumento en la variable independiente está asociado con una disminución en el log-odds de la variable dependiente.

10.3 Odds Ratios

Las odds ratios proporcionan una interpretación más intuitiva de los coeficientes de regresión logística. La odds ratio para una variable independiente X_j se calcula como e^{β_j} .

10.3.1 Cálculo de las Odds Ratios

$$OR_j = e^{\beta_j}$$

10.3.2 Interpretación de las Odds Ratios

- **OR > 1:** Un OR mayor que 1 indica que un aumento en la variable independiente está asociado con un aumento en las odds de la variable dependiente.
- **OR < 1:** Un OR menor que 1 indica que un aumento en la variable independiente está asociado con una disminución en las odds de la variable dependiente.
- **OR = 1:** Un OR igual a 1 indica que la variable independiente no tiene efecto sobre las odds de la variable dependiente.

10.4 Intervalos de Confianza

Los intervalos de confianza proporcionan una medida de la incertidumbre asociada con los estimadores de los coeficientes. Un intervalo de confianza del 95% para un coeficiente β_j indica que, en el 95% de las muestras, el intervalo contendrá el valor verdadero de β_j .

10.4.1 Cálculo de los Intervalos de Confianza

Para calcular un intervalo de confianza del 95% para un coeficiente β_j , utilizamos la fórmula:

$$\beta_j \pm 1.96 \cdot \text{SE}(\beta_j)$$

donde $\text{SE}(\beta_j)$ es el error estándar de β_j .

10.5 Significancia Estadística

La significancia estadística se utiliza para determinar si los coeficientes del modelo son significativamente diferentes de cero. Esto se evalúa mediante pruebas de hipótesis.

10.5.1 Prueba de Hipótesis

Para cada coeficiente β_j , la hipótesis nula H_0 es que $\beta_j = 0$. La hipótesis alternativa H_a es que $\beta_j \neq 0$.

10.5.2 P-valor

El p-valor indica la probabilidad de obtener un coeficiente tan extremo como el observado, asumiendo que la hipótesis nula es verdadera. Un p-valor menor que el nivel de significancia α (típicamente 0.05) indica que podemos rechazar la hipótesis nula.

10.6 Implementación en R

10.6.1 Cálculo de Coeficientes y Odds Ratios

```
# Cargar el paquete necesario
library(broom)
```

```
# Entrenar el modelo de regresión logística
model <- glm(var1 ~ ., data = dataTrain, family = "binomial")
```

```
# Coeficientes del modelo  
coef(model)
```

```
# Odds ratios  
exp(coef(model))
```

10.6.2 Intervalos de Confianza

```
# Intervalos de confianza para los coeficientes  
confint(model)
```

```
# Intervalos de confianza para las odds ratios  
exp(confint(model))
```

10.6.3 P-valores y Significancia Estadística

```
# Resumen del modelo con p-valores  
summary(model)
```

Chapter 11

Regresión Logística Multinomial y Análisis de Supervivencia

11.1 Introducción

La regresión logística multinomial y el análisis de supervivencia son extensiones de la regresión logística binaria. Este capítulo se enfoca en las técnicas y aplicaciones de estos métodos avanzados.

11.2 Regresión Logística Multinomial

La regresión logística multinomial se utiliza cuando la variable dependiente tiene más de dos categorías.

11.2.1 Modelo Multinomial

El modelo de regresión logística multinomial generaliza el modelo binario para manejar múltiples categorías. La probabilidad de que una observación pertenezca a la categoría k se expresa como:

$$P(Y = k) = \frac{e^{\beta_{0k} + \beta_{1k}X_1 + \dots + \beta_{nk}X_n}}{\sum_{j=1}^K e^{\beta_{0j} + \beta_{1j}X_1 + \dots + \beta_{nj}X_n}}$$

11.2.2 Estimación de Parámetros

Los coeficientes del modelo multinomial se estiman utilizando máxima verosimilitud, similar a la regresión logística binaria.

11.3 Análisis de Supervivencia

El análisis de supervivencia se utiliza para modelar el tiempo hasta que ocurre un evento de interés, como la muerte o la falla de un componente.

11.3.1 Función de Supervivencia

La función de supervivencia $S(t)$ describe la probabilidad de que una observación sobreviva más allá del tiempo t :

$$S(t) = P(T > t)$$

11.3.2 Modelo de Riesgos Proporcionales de Cox

El modelo de Cox es un modelo de regresión semiparamétrico utilizado para analizar datos de supervivencia:

$$h(t|X) = h_0(t)e^{\beta_1 X_1 + \dots + \beta_p X_p}$$

donde $h(t|X)$ es la tasa de riesgo en el tiempo t dado el vector de covariables X y $h_0(t)$ es la tasa de riesgo basal.

11.4 Implementación en R

11.4.1 Regresión Logística Multinomial

```
# Cargar el paquete necesario
library(nnet)

# Entrenar el modelo de regresión logística multinomial
model_multinom <- multinom(var1 ~ ., data = dataTrain)

# Resumen del modelo
summary(model_multinom)
```

11.4.2 Análisis de Supervivencia

```
# Cargar el paquete necesario
library(survival)

# Crear el objeto de supervivencia
surv_object <- Surv(time = data$time, event = data$status)

# Ajustar el modelo de Cox
model_cox <- coxph(surv_object ~ var1 + var2, data = data)

# Resumen del modelo
summary(model_cox)
```

Chapter 12

Implementación de Regresión Logística en Datos Reales

12.1 Introducción

Implementar un modelo de regresión logística en datos reales implica varias etapas, desde la limpieza de datos hasta la evaluación y validación del modelo. Este capítulo presenta un ejemplo práctico de la implementación de un modelo de regresión logística utilizando un conjunto de datos real.

12.2 Conjunto de Datos

Para este ejemplo, utilizaremos un conjunto de datos disponible públicamente que contiene información sobre clientes bancarios. El objetivo es predecir si un cliente suscribirá un depósito a plazo fijo.

12.3 Preparación de Datos

12.3.1 Carga y Exploración de Datos

Primero, cargamos y exploramos el conjunto de datos para entender su estructura y contenido.

```
# Cargar el paquete necesario
library(dplyr)

# Cargar el conjunto de datos
data <- read.csv("bank.csv")

# Explorar los datos
str(data)
summary(data)
```

12.3.2 Limpieza de Datos

El siguiente paso es limpiar los datos, lo que incluye tratar los valores faltantes y eliminar las duplicidades.

```
# Eliminar duplicados
data <- data %>% distinct()
```



```
# Imputar valores faltantes (si existen)
data <- data %>% mutate_if(is.numeric, ~ifelse(is.na(.), mean(., na.rm = TRUE), .))
```

12.3.3 Codificación de Variables Categóricas

Convertimos las variables categóricas en variables numéricas utilizando la codificación one-hot.

```
# Codificación one-hot de variables categóricas
data <- data %>% mutate(across(where(is.factor), ~ as.numeric(as.factor(.))))
```

12.4 División de Datos

Dividimos los datos en conjuntos de entrenamiento y prueba.

```
# Dividir los datos en conjuntos de entrenamiento y prueba
set.seed(123)
trainIndex <- createDataPartition(data$y, p = .8, list = FALSE, times = 1)
dataTrain <- data[trainIndex,]
dataTest <- data[-trainIndex,]
```

12.5 Entrenamiento del Modelo

Entrenamos un modelo de regresión logística utilizando el conjunto de entrenamiento.

```
# Entrenar el modelo de regresión logística
model <- glm(y ~ ., data = dataTrain, family = "binomial")

# Resumen del modelo
summary(model)
```

12.6 Evaluación del Modelo

Evaluamos el rendimiento del modelo utilizando el conjunto de prueba.

```
# Predicciones en el conjunto de prueba
predictions <- predict(model, dataTest, type = "response")

# Convertir probabilidades a etiquetas
predicted_labels <- ifelse(predictions > 0.5, 1, 0)

# Matriz de confusión
confusionMatrix(predicted_labels, dataTest$y)
```

12.7 Interpretación de los Resultados

Interpretamos los coeficientes del modelo y las odds ratios.

```
# Coeficientes del modelo  
coef(model)
```

```
# Odds ratios  
exp(coef(model))
```

Chapter 13

Resumen y Proyecto Final

13.1 Resumen de Conceptos Clave

En este curso, hemos cubierto una variedad de conceptos y técnicas esenciales para la regresión logística. Los conceptos clave incluyen:

- **Fundamentos de Probabilidad y Estadística:** Comprensión de distribuciones de probabilidad, medidas de tendencia central y dispersión, inferencia estadística y pruebas de hipótesis.
- **Regresión Logística:** Modelo de regresión logística binaria y multinomial, interpretación de coeficientes y odds ratios, métodos de estimación y validación.
- **Preparación de Datos:** Limpieza de datos, tratamiento de valores faltantes, codificación de variables categóricas y selección de variables.
- **Evaluación del Modelo:** Curva ROC, AUC, matriz de confusión, precisión, recall, F1-score y validación cruzada.
- **Diagnóstico del Modelo:** Análisis de residuos, influencia, multicolinealidad y ajuste de parámetros.
- **Análisis de Supervivencia:** Modelos de supervivencia, función de supervivencia y modelos de riesgos proporcionales de Cox.

13.2 Buenas Prácticas

Al implementar modelos de regresión logística, es importante seguir buenas prácticas para garantizar la precisión y la robustez de los modelos. Algunas buenas prácticas incluyen:

- **Exploración y Preparación de Datos:** Realizar un análisis exploratorio exhaustivo y preparar los datos adecuadamente antes de construir el modelo.
- **Evaluación y Validación del Modelo:** Utilizar métricas adecuadas para evaluar el rendimiento del modelo y validar el modelo utilizando técnicas como la validación cruzada.
- **Interpretación de Resultados:** Interpretar correctamente los coeficientes del modelo y las odds ratios, y comunicar los resultados de manera clara y concisa.
- **Revisión y Ajuste del Modelo:** Diagnosticar problemas en el modelo y ajustar los parámetros para mejorar el rendimiento.

13.3 Proyecto Final

Para aplicar los conceptos y técnicas aprendidos en este curso, te proponemos realizar un proyecto final utilizando un conjunto de datos de tu elección. El proyecto debe incluir las siguientes etapas:

13.3.1 Selección del Conjunto de Datos

Elige un conjunto de datos relevante que contenga una variable dependiente binaria o multinomial y varias variables independientes.

13.3.2 Exploración y Preparación de Datos

Realiza un análisis exploratorio de los datos y prepara los datos para el modelado. Esto incluye la limpieza de datos, el tratamiento de valores faltantes y la codificación de variables categóricas.

13.3.3 Entrenamiento y Evaluación del Modelo

Entrena un modelo de regresión logística utilizando el conjunto de datos preparado y evalúa su rendimiento utilizando métricas apropiadas.

13.3.4 Interpretación de Resultados

Interpreta los coeficientes del modelo y las odds ratios, y proporciona una explicación clara de los resultados.

13.3.5 Presentación del Proyecto

Presenta tu proyecto en un informe detallado que incluya la descripción del conjunto de datos, los pasos de preparación y modelado, los resultados del modelo y las conclusiones.

Part I

SEGUNDA PARTE: ANALISIS DE SUPERVIVENCIA

Chapter 14

Introducción al Análisis de Supervivencia

14.1 Conceptos Básicos

El análisis de supervivencia es una rama de la estadística que se ocupa del análisis del tiempo que transcurre hasta que ocurre un evento de interés, comúnmente referido como "tiempo de falla". Este campo es ampliamente utilizado en medicina, biología, ingeniería, ciencias sociales, y otros campos.

14.2 Definición de Eventos y Tiempos

En el análisis de supervivencia, un "evento" se refiere a la ocurrencia de un evento específico, como la muerte, la falla de un componente, la recaída de una enfermedad, etc. El "tiempo de supervivencia" es el tiempo que transcurre desde un punto de inicio definido hasta la ocurrencia del evento.

14.3 Censura

La censura ocurre cuando la información completa sobre el tiempo hasta el evento no está disponible para todos los individuos en el estudio. Hay tres tipos principales de censura:

- **Censura a la derecha:** Ocurre cuando el evento de interés no se ha observado para algunos sujetos antes del final del estudio.
- **Censura a la izquierda:** Ocurre cuando el evento de interés ocurrió antes del inicio del periodo de observación.
- **Censura por intervalo:** Ocurre cuando el evento de interés se sabe que ocurrió en un intervalo de tiempo, pero no se conoce el momento exacto.

14.4 Función de Supervivencia

La función de supervivencia, $S(t)$, se define como la probabilidad de que un individuo sobreviva más allá de un tiempo t . Matemáticamente, se expresa como:

$$S(t) = P(T > t)$$

donde T es una variable aleatoria que representa el tiempo hasta el evento. La función de supervivencia tiene las siguientes propiedades:

- $S(0) = 1$: Esto indica que al inicio (tiempo $t = 0$), la probabilidad de haber experimentado el evento es cero, por lo tanto, la supervivencia es del 100
- $\lim_{t \rightarrow \infty} S(t) = 0$: A medida que el tiempo tiende al infinito, la probabilidad de que cualquier individuo aún no haya experimentado el evento tiende a cero.
- $S(t)$ es una función no creciente: Esto significa que a medida que el tiempo avanza, la probabilidad de supervivencia no aumenta.

14.5 Función de Densidad de Probabilidad

La función de densidad de probabilidad $f(t)$ describe la probabilidad de que el evento ocurra en un instante de tiempo específico. Se define como:

$$f(t) = \frac{dF(t)}{dt}$$

donde $F(t)$ es la función de distribución acumulada, $F(t) = P(T \leq t)$. La relación entre $S(t)$ y $f(t)$ es:

$$f(t) = -\frac{dS(t)}{dt}$$

14.6 Función de Riesgo

La función de riesgo, $\lambda(t)$, también conocida como función de tasa de fallas o hazard rate, se define como la tasa instantánea de ocurrencia del evento en el tiempo t , dado que el individuo ha sobrevivido hasta el tiempo t . Matemáticamente, se expresa como:

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t}$$

Esto se puede reescribir usando $f(t)$ y $S(t)$ como:

$$\lambda(t) = \frac{f(t)}{S(t)}$$

14.7 Relación entre Función de Supervivencia y Función de Riesgo

La función de supervivencia y la función de riesgo están relacionadas a través de la siguiente ecuación:

$$S(t) = \exp \left(- \int_0^t \lambda(u) du \right)$$

Esta fórmula se deriva del hecho de que la función de supervivencia es la probabilidad acumulativa de no haber experimentado el evento hasta el tiempo t , y $\lambda(t)$ es la tasa instantánea de ocurrencia del evento.

La función de riesgo también puede ser expresada como:

$$\lambda(t) = -\frac{d}{dt} \log S(t)$$

14.8 Deducción de la Función de Supervivencia

La relación entre la función de supervivencia y la función de riesgo se puede deducir integrando la función de riesgo:

$$\begin{aligned}S(t) &= \exp\left(-\int_0^t \lambda(u) du\right) \\ \log S(t) &= -\int_0^t \lambda(u) du \\ \frac{d}{dt} \log S(t) &= -\lambda(t) \\ \lambda(t) &= -\frac{d}{dt} \log S(t)\end{aligned}$$

14.9 Ejemplo de Cálculo

Supongamos que tenemos una muestra de tiempos de supervivencia T_1, T_2, \dots, T_n . Podemos estimar la función de supervivencia empírica como:

$$\hat{S}(t) = \frac{\text{Número de individuos que sobreviven más allá de } t}{\text{Número total de individuos en riesgo en } t}$$

y la función de riesgo empírica como:

$$\hat{\lambda}(t) = \frac{\text{Número de eventos en } t}{\text{Número de individuos en riesgo en } t}$$

14.10 Conclusión

El análisis de supervivencia es una herramienta poderosa para analizar datos de tiempo hasta evento. Entender los conceptos básicos como la función de supervivencia y la función de riesgo es fundamental para el análisis más avanzado.

Chapter 15

Función de Supervivencia y Función de Riesgo

15.1 Introducción

Este capítulo profundiza en la definición y propiedades de la función de supervivencia y la función de riesgo, dos conceptos fundamentales en el análisis de supervivencia. Entender estas funciones y su relación es crucial para modelar y analizar datos de tiempo hasta evento.

15.2 Función de Supervivencia

La función de supervivencia, $S(t)$, describe la probabilidad de que un individuo sobreviva más allá de un tiempo t . Formalmente, se define como:

$$S(t) = P(T > t)$$

donde T es una variable aleatoria que representa el tiempo hasta el evento.

15.2.1 Propiedades de la Función de Supervivencia

La función de supervivencia tiene varias propiedades importantes:

- $S(0) = 1$: Indica que la probabilidad de haber experimentado el evento en el tiempo 0 es cero.
- $\lim_{t \rightarrow \infty} S(t) = 0$: A medida que el tiempo tiende al infinito, la probabilidad de supervivencia tiende a cero.
- $S(t)$ es una función no creciente: A medida que el tiempo avanza, la probabilidad de supervivencia no aumenta.

15.2.2 Derivación de $S(t)$

Si la función de densidad de probabilidad $f(t)$ del tiempo de supervivencia T es conocida, la función de supervivencia puede derivarse como:

$$\begin{aligned} S(t) &= P(T > t) \\ &= 1 - P(T \leq t) \\ &= 1 - F(t) \\ &= 1 - \int_0^t f(u) du \end{aligned}$$

donde $F(t)$ es la función de distribución acumulada.

15.2.3 Ejemplo de Cálculo de $S(t)$

Consideremos un ejemplo donde el tiempo de supervivencia T sigue una distribución exponencial con tasa λ . La función de densidad de probabilidad $f(t)$ es:

$$f(t) = \lambda e^{-\lambda t}, \quad t \geq 0$$

La función de distribución acumulada $F(t)$ es:

$$F(t) = \int_0^t \lambda e^{-\lambda u} du = 1 - e^{-\lambda t}$$

Por lo tanto, la función de supervivencia $S(t)$ es:

$$S(t) = 1 - F(t) = e^{-\lambda t}$$

15.3 Función de Riesgo

La función de riesgo, $\lambda(t)$, proporciona la tasa instantánea de ocurrencia del evento en el tiempo t , dado que el individuo ha sobrevivido hasta el tiempo t . Matemáticamente, se define como:

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t}$$

15.3.1 Relación entre $\lambda(t)$ y $f(t)$

La función de riesgo se puede relacionar con la función de densidad de probabilidad $f(t)$ y la función de supervivencia $S(t)$ de la siguiente manera:

$$\lambda(t) = \frac{f(t)}{S(t)}$$

15.3.2 Derivación de $\lambda(t)$

La derivación de $\lambda(t)$ se basa en la definición condicional de la probabilidad:

$$\begin{aligned} \lambda(t) &= \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t} \\ &= \lim_{\Delta t \rightarrow 0} \frac{\frac{P(t \leq T < t + \Delta t \text{ y } T \geq t)}{P(T \geq t)}}{\Delta t} \\ &= \lim_{\Delta t \rightarrow 0} \frac{\frac{P(t \leq T < t + \Delta t)}{P(T \geq t)}}{\Delta t} \\ &= \frac{f(t)}{S(t)} \end{aligned}$$

15.4 Relación entre Función de Supervivencia y Función de Riesgo

La función de supervivencia y la función de riesgo están estrechamente relacionadas. La relación se expresa mediante la siguiente ecuación:

$$S(t) = \exp \left(- \int_0^t \lambda(u) du \right)$$

15.4.1 Deducción de la Relación

Para deducir esta relación, consideramos la derivada logarítmica de la función de supervivencia:

$$\begin{aligned} S(t) &= \exp \left(- \int_0^t \lambda(u) du \right) \\ \log S(t) &= - \int_0^t \lambda(u) du \\ \frac{d}{dt} \log S(t) &= -\lambda(t) \\ \lambda(t) &= -\frac{d}{dt} \log S(t) \end{aligned}$$

15.5 Interpretación de la Función de Riesgo

La función de riesgo, $\lambda(t)$, se interpreta como la tasa instantánea de ocurrencia del evento por unidad de tiempo, dado que el individuo ha sobrevivido hasta el tiempo t . Es una medida local del riesgo de falla en un instante específico.

15.5.1 Ejemplo de Cálculo de $\lambda(t)$

Consideremos nuevamente el caso donde el tiempo de supervivencia T sigue una distribución exponencial con tasa λ . La función de densidad de probabilidad $f(t)$ es:

$$f(t) = \lambda e^{-\lambda t}$$

La función de supervivencia $S(t)$ es:

$$S(t) = e^{-\lambda t}$$

La función de riesgo $\lambda(t)$ se calcula como:

$$\begin{aligned} \lambda(t) &= \frac{f(t)}{S(t)} \\ &= \frac{\lambda e^{-\lambda t}}{e^{-\lambda t}} \\ &= \lambda \end{aligned}$$

En este caso, $\lambda(t)$ es constante y igual a λ , lo que es una característica de la distribución exponencial.

15.6 Funciones de Riesgo Acumulada y Media Residual

La función de riesgo acumulada $H(t)$ se define como:

$$H(t) = \int_0^t \lambda(u) du$$

Esta función proporciona la suma acumulada de la tasa de riesgo hasta el tiempo t .

La función de vida media residual $e(t)$ se define como la esperanza del tiempo de vida restante dado que el individuo ha sobrevivido hasta el tiempo t :

$$e(t) = \mathbb{E}[T - t \mid T > t] = \int_t^\infty S(u) du$$

15.7 Ejemplo de Cálculo de Función de Riesgo Acumulada y Vida Media Residual

Consideremos nuevamente la distribución exponencial con tasa λ . La función de riesgo acumulada $H(t)$ es:

$$\begin{aligned} H(t) &= \int_0^t \lambda du \\ &= \lambda t \end{aligned}$$

La función de vida media residual $e(t)$ es:

$$\begin{aligned} e(t) &= \int_t^\infty e^{-\lambda u} du \\ &= \left[\frac{-1}{\lambda} e^{-\lambda u} \right]_t^\infty \\ &= \frac{1}{\lambda} e^{-\lambda t} \\ &= \frac{1}{\lambda} \end{aligned}$$

En este caso, la vida media residual es constante e igual a $\frac{1}{\lambda}$, otra característica de la distribución exponencial.

15.8 Conclusión

La función de supervivencia y la función de riesgo son herramientas fundamentales en el análisis de supervivencia. Entender su definición, propiedades, y la relación entre ellas es esencial para modelar y analizar correctamente los datos de tiempo hasta evento. Las funciones de riesgo acumulada y vida media residual proporcionan información adicional sobre la dinámica del riesgo a lo largo del tiempo.

Chapter 16

Estimador de Kaplan-Meier

16.1 Introducción

El estimador de Kaplan-Meier, también conocido como la función de supervivencia empírica, es una herramienta no paramétrica para estimar la función de supervivencia a partir de datos censurados. Este método es especialmente útil cuando los tiempos de evento están censurados a la derecha.

16.2 Definición del Estimador de Kaplan-Meier

El estimador de Kaplan-Meier se define como:

$$\hat{S}(t) = \prod_{t_i \leq t} \left(1 - \frac{d_i}{n_i}\right)$$

donde:

- t_i es el tiempo del i -ésimo evento,
- d_i es el número de eventos que ocurren en t_i ,
- n_i es el número de individuos en riesgo justo antes de t_i .

16.3 Propiedades del Estimador de Kaplan-Meier

El estimador de Kaplan-Meier tiene las siguientes propiedades:

- Es una función escalonada que disminuye en los tiempos de los eventos observados.
- Puede manejar datos censurados a la derecha.
- Proporciona una estimación no paramétrica de la función de supervivencia.

16.3.1 Función Escalonada

La función escalonada del estimador de Kaplan-Meier significa que $\hat{S}(t)$ permanece constante entre los tiempos de los eventos y disminuye en los tiempos de los eventos. Matemáticamente, si t_i es el tiempo del i -ésimo evento, entonces:

$$\hat{S}(t) = \hat{S}(t_i) \quad \text{para } t_i \leq t < t_{i+1}$$

16.3.2 Manejo de Datos Censurados

El estimador de Kaplan-Meier maneja datos censurados a la derecha al ajustar la estimación de la función de supervivencia sólo en los tiempos en que ocurren eventos. Si un individuo es censurado antes de experimentar el evento, no contribuye a la disminución de $\hat{S}(t)$ en el tiempo de censura. Esto asegura que la censura no sesga la estimación de la supervivencia.

16.3.3 Estimación No Paramétrica

El estimador de Kaplan-Meier es no paramétrico porque no asume ninguna forma específica para la distribución de los tiempos de supervivencia. En cambio, utiliza la información empírica disponible para estimar la función de supervivencia.

16.4 Deducción del Estimador de Kaplan-Meier

La deducción del estimador de Kaplan-Meier se basa en el principio de probabilidad condicional. Consideremos un conjunto de tiempos de supervivencia observados t_1, t_2, \dots, t_k con eventos en cada uno de estos tiempos. El estimador de la probabilidad de supervivencia más allá del tiempo t es el producto de las probabilidades de sobrevivir más allá de cada uno de los tiempos de evento observados hasta t .

16.4.1 Probabilidad Condicional

La probabilidad de sobrevivir más allá de t_i , dado que el individuo ha sobrevivido justo antes de t_i , es:

$$P(T > t_i \mid T \geq t_i) = 1 - \frac{d_i}{n_i}$$

donde d_i es el número de eventos en t_i y n_i es el número de individuos en riesgo justo antes de t_i .

16.4.2 Producto de Probabilidades Condicionales

La probabilidad de sobrevivir más allá de un tiempo t cualquiera, dada la secuencia de tiempos de evento, es el producto de las probabilidades condicionales de sobrevivir más allá de cada uno de los tiempos de evento observados hasta t . Así, el estimador de Kaplan-Meier se obtiene como:

$$\hat{S}(t) = \prod_{t_i \leq t} \left(1 - \frac{d_i}{n_i}\right)$$

16.5 Ejemplo de Cálculo

Supongamos que tenemos los siguientes tiempos de supervivencia observados para cinco individuos: 2, 3, 5, 7, 8. Supongamos además que tenemos censura a la derecha en el tiempo 10. Los tiempos de evento y el número de individuos en riesgo justo antes de cada evento son:

Tiempo (t_i)	Eventos (d_i)	En Riesgo (n_i)
2	1	5
3	1	4
5	1	3
7	1	2
8	1	1

Table 16.1: Ejemplo de cálculo del estimador de Kaplan-Meier

Usando estos datos, el estimador de Kaplan-Meier se calcula como:

$$\begin{aligned}
\hat{S}(2) &= 1 - \frac{1}{5} = 0.8 \\
\hat{S}(3) &= 0.8 \times \left(1 - \frac{1}{4}\right) = 0.8 \times 0.75 = 0.6 \\
\hat{S}(5) &= 0.6 \times \left(1 - \frac{1}{3}\right) = 0.6 \times 0.6667 = 0.4 \\
\hat{S}(7) &= 0.4 \times \left(1 - \frac{1}{2}\right) = 0.4 \times 0.5 = 0.2 \\
\hat{S}(8) &= 0.2 \times \left(1 - \frac{1}{1}\right) = 0.2 \times 0 = 0
\end{aligned}$$

16.6 Intervalos de Confianza para el Estimador de Kaplan-Meier

Para calcular intervalos de confianza para el estimador de Kaplan-Meier, se puede usar la transformación logarítmica y la aproximación normal. Un intervalo de confianza aproximado para $\log(-\log(\hat{S}(t)))$ se obtiene como:

$$\log(-\log(\hat{S}(t))) \pm z_{\alpha/2} \sqrt{\frac{1}{d_i(n_i - d_i)}}$$

donde $z_{\alpha/2}$ es el percentil correspondiente de la distribución normal estándar.

16.7 Transformación Logarítmica Inversa

La transformación logarítmica inversa se utiliza para obtener los límites del intervalo de confianza para $S(t)$:

$$\hat{S}(t) = \exp \left(- \exp \left(\log(-\log(\hat{S}(t))) \pm z_{\alpha/2} \sqrt{\frac{1}{d_i(n_i - d_i)}} \right) \right)$$

16.8 Cálculo Detallado de Intervalos de Confianza

Para un cálculo más detallado de los intervalos de confianza, consideremos un tiempo específico t_j . La varianza del estimador de Kaplan-Meier en t_j se puede estimar usando Greenwood's formula:

$$\text{Var}(\hat{S}(t_j)) = \hat{S}(t_j)^2 \sum_{t_i \leq t_j} \frac{d_i}{n_i(n_i - d_i)}$$

El intervalo de confianza aproximado para $\hat{S}(t_j)$ es entonces:

$$\hat{S}(t_j) \pm z_{\alpha/2} \sqrt{\text{Var}(\hat{S}(t_j))}$$

16.9 Ejemplo de Intervalo de Confianza

Supongamos que en el ejemplo anterior queremos calcular el intervalo de confianza para $\hat{S}(3)$. Primero, calculamos la varianza:

$$\begin{aligned} \text{Var}(\hat{S}(3)) &= \hat{S}(3)^2 \left(\frac{1}{5 \times 4} + \frac{1}{4 \times 3} \right) \\ &= 0.6^2 \left(\frac{1}{20} + \frac{1}{12} \right) \\ &= 0.36 (0.05 + 0.0833) \\ &= 0.36 \times 0.1333 \\ &= 0.048 \end{aligned}$$

El intervalo de confianza es entonces:

$$0.6 \pm 1.96\sqrt{0.048} = 0.6 \pm 1.96 \times 0.219 = 0.6 \pm 0.429$$

Por lo tanto, el intervalo de confianza para $\hat{S}(3)$ es aproximadamente $(0.171, 1.029)$. Dado que una probabilidad no puede exceder 1, ajustamos el intervalo a $(0.171, 1.0)$.

16.10 Interpretación del Estimador de Kaplan-Meier

El estimador de Kaplan-Meier proporciona una estimación empírica de la función de supervivencia que es fácil de interpretar y calcular. Su capacidad para manejar datos censurados lo hace especialmente útil en estudios de supervivencia.

16.11 Conclusión

El estimador de Kaplan-Meier es una herramienta poderosa para estimar la función de supervivencia en presencia de datos censurados. Su cálculo es relativamente sencillo y proporciona una estimación no paramétrica robusta de la supervivencia a lo largo del tiempo. La interpretación adecuada de este estimador y su intervalo de confianza asociado es fundamental para el análisis de datos de supervivencia.

Chapter 17

Comparación de Curvas de Supervivencia

17.1 Introducción

Comparar curvas de supervivencia es crucial para determinar si existen diferencias significativas en las tasas de supervivencia entre diferentes grupos. Las pruebas de hipótesis, como el test de log-rank, son herramientas comunes para esta comparación.

17.2 Test de Log-rank

El test de log-rank se utiliza para comparar las curvas de supervivencia de dos o más grupos. La hipótesis nula es que no hay diferencia en las funciones de riesgo entre los grupos.

17.2.1 Fórmula del Test de Log-rank

El estadístico del test de log-rank se define como:

$$\chi^2 = \frac{\left(\sum_{i=1}^k (O_i - E_i)\right)^2}{\sum_{i=1}^k V_i}$$

donde:

- O_i es el número observado de eventos en el grupo i .
- E_i es el número esperado de eventos en el grupo i .
- V_i es la varianza del número de eventos en el grupo i .

17.2.2 Cálculo de E_i y V_i

El número esperado de eventos E_i y la varianza V_i se calculan como:

$$E_i = \frac{d_i \cdot n_i}{n}$$
$$V_i = \frac{d_i \cdot (n - d_i) \cdot n_i \cdot (n - n_i)}{n^2 \cdot (n - 1)}$$

donde:

- d_i es el número total de eventos en el grupo i .
- n_i es el número de individuos en riesgo en el grupo i .
- n es el número total de individuos en todos los grupos.

17.3 Ejemplo de Cálculo del Test de Log-rank

Supongamos que tenemos dos grupos con los siguientes datos de eventos:

Grupo	Tiempo (t_i)	Eventos (O_i)	En Riesgo (n_i)
1	2	1	5
1	4	1	4
2	3	1	4
2	5	1	3

Table 17.1: Ejemplo de datos para el test de log-rank

Calculemos E_i y V_i para cada grupo:

$$\begin{aligned}
 E_1 &= \frac{2 \cdot 5}{9} + \frac{2 \cdot 4}{8} = \frac{10}{9} + \frac{8}{8} = 1.11 + 1 = 2.11 \\
 V_1 &= \frac{2 \cdot 7 \cdot 5 \cdot 4}{81 \cdot 8} = \frac{2 \cdot 7 \cdot 5 \cdot 4}{648} = \frac{280}{648} = 0.432 \\
 E_2 &= \frac{2 \cdot 4}{9} + \frac{2 \cdot 3}{8} = \frac{8}{9} + \frac{6}{8} = 0.89 + 0.75 = 1.64 \\
 V_2 &= \frac{2 \cdot 7 \cdot 4 \cdot 4}{81 \cdot 8} = \frac{2 \cdot 7 \cdot 4 \cdot 4}{648} = \frac{224}{648} = 0.346
 \end{aligned}$$

El estadístico de log-rank se calcula como:

$$\begin{aligned}
 \chi^2 &= \frac{((1 - 2.11) + (1 - 1.64))^2}{0.432 + 0.346} \\
 &= \frac{(-1.11 - 0.64)^2}{0.778} \\
 &= \frac{3.04}{0.778} \\
 &= 3.91
 \end{aligned}$$

El valor p se puede obtener comparando χ^2 con una distribución χ^2 con un grado de libertad (dado que estamos comparando dos grupos).

17.4 Interpretación del Test de Log-rank

Un valor p pequeño (generalmente menos de 0.05) indica que hay una diferencia significativa en las curvas de supervivencia entre los grupos. Un valor p grande sugiere que no hay suficiente evidencia para rechazar la hipótesis nula de que las curvas de supervivencia son iguales.

17.5 Pruebas Alternativas

Además del test de log-rank, existen otras pruebas para comparar curvas de supervivencia, como el test de Wilcoxon (Breslow), que da más peso a los eventos en tiempos tempranos.

17.6 Conclusión

El test de log-rank es una herramienta esencial para comparar curvas de supervivencia entre diferentes grupos. Su cálculo se basa en la diferencia entre los eventos observados y esperados en cada grupo, y su interpretación puede ayudar a identificar diferencias significativas en la supervivencia.

Chapter 18

Modelos de Riesgos Proporcionales de Cox

18.1 Introducción

El modelo de riesgos proporcionales de Cox, propuesto por David Cox en 1972, es una de las herramientas más utilizadas en el análisis de supervivencia. Este modelo permite evaluar el efecto de varias covariables en el tiempo hasta el evento, sin asumir una forma específica para la distribución de los tiempos de supervivencia.

18.2 Definición del Modelo de Cox

El modelo de Cox se define como:

$$\lambda(t | X) = \lambda_0(t) \exp(\beta^T X)$$

donde:

- $\lambda(t | X)$ es la función de riesgo en el tiempo t dado el vector de covariables X .
- $\lambda_0(t)$ es la función de riesgo basal en el tiempo t .
- β es el vector de coeficientes del modelo.
- X es el vector de covariables.

18.3 Supuesto de Proporcionalidad de Riesgos

El modelo de Cox asume que las razones de riesgo entre dos individuos son constantes a lo largo del tiempo. Matemáticamente, si X_i y X_j son las covariables de dos individuos, la razón de riesgos se expresa como:

$$\frac{\lambda(t | X_i)}{\lambda(t | X_j)} = \frac{\lambda_0(t) \exp(\beta^T X_i)}{\lambda_0(t) \exp(\beta^T X_j)} = \exp(\beta^T (X_i - X_j))$$

18.4 Estimación de los Parámetros

Los parámetros β se estiman utilizando el método de máxima verosimilitud parcial. La función de verosimilitud parcial se define como:

$$L(\beta) = \prod_{i=1}^k \frac{\exp(\beta^T X_i)}{\sum_{j \in R(t_i)} \exp(\beta^T X_j)}$$

donde $R(t_i)$ es el conjunto de individuos en riesgo en el tiempo t_i .

18.4.1 Función de Log-Verosimilitud Parcial

La función de log-verosimilitud parcial es:

$$\log L(\beta) = \sum_{i=1}^k \left(\beta^T X_i - \log \sum_{j \in R(t_i)} \exp(\beta^T X_j) \right)$$

18.4.2 Derivadas Parciales y Maximización

Para encontrar los estimadores de máxima verosimilitud, resolvemos el sistema de ecuaciones obtenido al igualar a cero las derivadas parciales de $\log L(\beta)$ con respecto a β :

$$\frac{\partial \log L(\beta)}{\partial \beta} = \sum_{i=1}^k \left(X_i - \frac{\sum_{j \in R(t_i)} X_j \exp(\beta^T X_j)}{\sum_{j \in R(t_i)} \exp(\beta^T X_j)} \right) = 0$$

18.5 Interpretación de los Coeficientes

Cada coeficiente β_i representa el logaritmo de la razón de riesgos asociado con un incremento unitario en la covariable X_i . Un valor positivo de β_i indica que un aumento en X_i incrementa el riesgo del evento, mientras que un valor negativo indica una reducción del riesgo.

18.6 Evaluación del Modelo

El modelo de Cox se evalúa utilizando varias técnicas, como el análisis de residuos de Schoenfeld para verificar el supuesto de proporcionalidad de riesgos, y el uso de curvas de supervivencia estimadas para evaluar la bondad de ajuste.

18.6.1 Residuos de Schoenfeld

Los residuos de Schoenfeld se utilizan para evaluar la proporcionalidad de riesgos. Para cada evento en el tiempo t_i , el residuo de Schoenfeld para la covariable X_j se define como:

$$r_{ij} = X_{ij} - \hat{X}_{ij}$$

donde \hat{X}_{ij} es la covariable ajustada.

18.6.2 Curvas de Supervivencia Ajustadas

Las curvas de supervivencia ajustadas se obtienen utilizando la función de riesgo basal estimada y los coeficientes del modelo. La función de supervivencia ajustada se define como:

$$\hat{S}(t | X) = \hat{S}_0(t)^{\exp(\beta^T X)}$$

donde $\hat{S}_0(t)$ es la función de supervivencia basal estimada.

18.7 Ejemplo de Aplicación del Modelo de Cox

Consideremos un ejemplo con tres covariables: edad, sexo y tratamiento. Supongamos que los datos se ajustan a un modelo de Cox y obtenemos los siguientes coeficientes:

$$\hat{\beta}_{edad} = 0.02, \quad \hat{\beta}_{sexo} = -0.5, \quad \hat{\beta}_{tratamiento} = 1.2$$

La función de riesgo ajustada se expresa como:

$$\lambda(t | X) = \lambda_0(t) \exp(0.02 \cdot \text{edad} - 0.5 \cdot \text{sexo} + 1.2 \cdot \text{tratamiento})$$

18.8 Conclusión

El modelo de riesgos proporcionales de Cox es una herramienta poderosa para analizar datos de supervivencia con múltiples covariables. Su flexibilidad y la falta de suposiciones fuertes sobre la distribución de los tiempos de supervivencia lo hacen ampliamente aplicable en diversas disciplinas.

Chapter 19

Diagnóstico y Validación de Modelos de Cox

19.1 Introducción

Una vez ajustado un modelo de Cox, es crucial realizar diagnósticos y validaciones para asegurar que el modelo es apropiado y que los supuestos subyacentes son válidos. Esto incluye la verificación del supuesto de proporcionalidad de riesgos y la evaluación del ajuste del modelo.

19.2 Supuesto de Proporcionalidad de Riesgos

El supuesto de proporcionalidad de riesgos implica que la razón de riesgos entre dos individuos es constante a lo largo del tiempo. Si este supuesto no se cumple, las inferencias hechas a partir del modelo pueden ser incorrectas.

19.2.1 Residuos de Schoenfeld

Los residuos de Schoenfeld se utilizan para evaluar la proporcionalidad de riesgos. Para cada evento en el tiempo t_i , el residuo de Schoenfeld para la covariable X_j se define como:

$$r_{ij} = X_{ij} - \hat{X}_{ij}$$

donde \hat{X}_{ij} es la covariable ajustada. Si los residuos de Schoenfeld no muestran una tendencia sistemática cuando se trazan contra el tiempo, el supuesto de proporcionalidad de riesgos es razonable.

19.3 Bondad de Ajuste

La bondad de ajuste del modelo de Cox se evalúa comparando las curvas de supervivencia observadas y ajustadas, y utilizando estadísticas de ajuste global.

19.3.1 Curvas de Supervivencia Ajustadas

Las curvas de supervivencia ajustadas se obtienen utilizando la función de riesgo basal estimada y los coeficientes del modelo. La función de supervivencia ajustada se define como:

$$\hat{S}(t | X) = \hat{S}_0(t)^{\exp(\beta^T X)}$$

donde $\hat{S}_0(t)$ es la función de supervivencia basal estimada. Comparar estas curvas con las curvas de Kaplan-Meier para diferentes niveles de las covariables puede proporcionar una validación visual del ajuste del modelo.

19.3.2 Estadísticas de Ajuste Global

Las estadísticas de ajuste global, como el test de la desviación y el test de la bondad de ajuste de Grambsch y Therneau, se utilizan para evaluar el ajuste global del modelo de Cox.

19.4 Diagnóstico de Influencia

El diagnóstico de influencia identifica observaciones individuales que tienen un gran impacto en los estimados del modelo. Los residuos de devianza y los residuos de martingala se utilizan comúnmente para este propósito.

19.4.1 Residuos de Deviance

Los residuos de deviance se definen como:

$$D_i = \text{sign}(O_i - E_i) \sqrt{-2 \left(O_i \log \frac{O_i}{E_i} - (O_i - E_i) \right)}$$

donde O_i es el número observado de eventos y E_i es el número esperado de eventos. Observaciones con residuos de deviance grandes en valor absoluto pueden ser influyentes.

19.4.2 Residuos de Martingala

Los residuos de martingala se definen como:

$$M_i = O_i - E_i$$

donde O_i es el número observado de eventos y E_i es el número esperado de eventos. Los residuos de martingala se utilizan para detectar observaciones que no se ajustan bien al modelo.

19.5 Ejemplo de Diagnóstico

Consideremos un modelo de Cox ajustado con las covariables edad, sexo y tratamiento. Para diagnosticar la influencia de observaciones individuales, calculamos los residuos de deviance y martingala para cada observación.

Observación	Edad	Sexo	Tratamiento	Residuo de Deviance
1	50	0	1	1.2
2	60	1	0	-0.5
3	45	0	1	-1.8
4	70	1	0	0.3

Table 19.1: Residuos de deviance para observaciones individuales

Observaciones con residuos de deviance grandes en valor absoluto (como la observación 3) pueden ser influyentes y requieren una revisión adicional.

19.6 Conclusión

El diagnóstico y la validación son pasos críticos en el análisis de modelos de Cox. Evaluar el supuesto de proporcionalidad de riesgos, la bondad de ajuste y la influencia de observaciones individuales asegura que las inferencias y conclusiones derivadas del modelo sean válidas y fiables.

Chapter 20

Modelos Acelerados de Fallos

20.1 Introducción

Los modelos acelerados de fallos (AFT) son una alternativa a los modelos de riesgos proporcionales de Cox. En lugar de asumir que las covariables afectan la tasa de riesgo, los modelos AFT asumen que las covariables multiplican el tiempo de supervivencia por una constante.

20.2 Definición del Modelo AFT

Un modelo AFT se expresa como:

$$T = T_0 \exp(\beta^T X)$$

donde:

- T es el tiempo de supervivencia observado.
- T_0 es el tiempo de supervivencia bajo condiciones basales.
- β es el vector de coeficientes del modelo.
- X es el vector de covariables.

20.2.1 Transformación Logarítmica

El modelo AFT se puede transformar logarítmicamente para obtener una forma lineal:

$$\log(T) = \log(T_0) + \beta^T X$$

20.3 Estimación de los Parámetros

Los parámetros del modelo AFT se estiman utilizando el método de máxima verosimilitud. La función de verosimilitud se define como:

$$L(\beta) = \prod_{i=1}^n f(t_i | X_i; \beta)$$

donde $f(t_i | X_i; \beta)$ es la función de densidad de probabilidad del tiempo de supervivencia t_i dado el vector de covariables X_i y los parámetros β .

20.3.1 Función de Log-Verosimilitud

La función de log-verosimilitud es:

$$\log L(\beta) = \sum_{i=1}^n \log f(t_i | X_i; \beta)$$

20.3.2 Maximización de la Verosimilitud

Los estimadores de máxima verosimilitud se obtienen resolviendo el sistema de ecuaciones obtenido al igualar a cero las derivadas parciales de $\log L(\beta)$ con respecto a β :

$$\frac{\partial \log L(\beta)}{\partial \beta} = 0$$

20.4 Distribuciones Comunes en Modelos AFT

En los modelos AFT, el tiempo de supervivencia T puede seguir varias distribuciones comunes, como la exponencial, Weibull, log-normal y log-logística. Cada una de estas distribuciones tiene diferentes propiedades y aplicaciones.

20.4.1 Modelo Exponencial AFT

En un modelo exponencial AFT, el tiempo de supervivencia T sigue una distribución exponencial con parámetro λ :

$$f(t) = \lambda \exp(-\lambda t)$$

La función de supervivencia es:

$$S(t) = \exp(-\lambda t)$$

La transformación logarítmica del tiempo de supervivencia es:

$$\log(T) = \log\left(\frac{1}{\lambda}\right) + \beta^T X$$

20.4.2 Modelo Weibull AFT

En un modelo Weibull AFT, el tiempo de supervivencia T sigue una distribución Weibull con parámetros λ y k :

$$f(t) = \lambda k t^{k-1} \exp(-\lambda t^k)$$

La función de supervivencia es:

$$S(t) = \exp(-\lambda t^k)$$

La transformación logarítmica del tiempo de supervivencia es:

$$\log(T) = \log\left(\left(\frac{1}{\lambda}\right)^{1/k}\right) + \frac{\beta^T X}{k}$$

20.5 Interpretación de los Coeficientes

En los modelos AFT, los coeficientes β_i se interpretan como factores multiplicativos del tiempo de supervivencia. Un valor positivo de β_i indica que un aumento en la covariable X_i incrementa el tiempo de supervivencia, mientras que un valor negativo indica una reducción del tiempo de supervivencia.

20.6 Ejemplo de Aplicación del Modelo AFT

Consideremos un ejemplo con tres covariables: edad, sexo y tratamiento. Supongamos que los datos se ajustan a un modelo Weibull AFT y obtenemos los siguientes coeficientes:

$$\hat{\beta}_{edad} = -0.02, \quad \hat{\beta}_{sexo} = 0.5, \quad \hat{\beta}_{tratamiento} = -1.2$$

La función de supervivencia ajustada se expresa como:

$$S(t | X) = \exp \left(- \left(\frac{t \exp(-0.02 \cdot edad + 0.5 \cdot sexo - 1.2 \cdot tratamiento)}{\lambda} \right)^k \right)$$

20.7 Conclusión

Los modelos AFT proporcionan una alternativa flexible a los modelos de riesgos proporcionales de Cox. Su enfoque en la multiplicación del tiempo de supervivencia por una constante permite una interpretación intuitiva y aplicaciones en diversas áreas.

Chapter 21

Análisis Multivariado de Supervivencia

21.1 Introducción

El análisis multivariado de supervivencia extiende los modelos de supervivencia para incluir múltiples covariables, permitiendo evaluar su efecto simultáneo sobre el tiempo hasta el evento. Los modelos de Cox y AFT son comúnmente utilizados en este contexto.

21.2 Modelo de Cox Multivariado

El modelo de Cox multivariado se define como:

$$\lambda(t | X) = \lambda_0(t) \exp(\beta^T X)$$

donde X es un vector de covariables.

21.2.1 Estimación de los Parámetros

Los parámetros β se estiman utilizando el método de máxima verosimilitud parcial, como se discutió anteriormente. La función de verosimilitud parcial se maximiza para obtener los estimadores de los coeficientes.

21.3 Modelo AFT Multivariado

El modelo AFT multivariado se expresa como:

$$T = T_0 \exp(\beta^T X)$$

21.3.1 Estimación de los Parámetros

Los parámetros β se estiman utilizando el método de máxima verosimilitud, similar al caso univariado. La función de verosimilitud se maximiza para obtener los estimadores de los coeficientes.

21.4 Interacción y Efectos No Lineales

En el análisis multivariado, es importante considerar la posibilidad de interacciones entre covariables y efectos no lineales. Estos se pueden incluir en los modelos extendiendo las funciones de riesgo o supervivencia.

21.4.1 Interacciones

Las interacciones entre covariables se pueden modelar añadiendo términos de interacción en el modelo:

$$\lambda(t | X) = \lambda_0(t) \exp(\beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2)$$

donde $X_1 X_2$ es el término de interacción.

21.4.2 Efectos No Lineales

Los efectos no lineales se pueden modelar utilizando funciones no lineales de las covariables, como polinomios o splines:

$$\lambda(t | X) = \lambda_0(t) \exp(\beta_1 X + \beta_2 X^2)$$

21.5 Selección de Variables

La selección de variables es crucial en el análisis multivariado para evitar el sobreajuste y mejorar la interpretabilidad del modelo. Métodos como la regresión hacia atrás, la regresión hacia adelante y la selección por criterios de información (AIC, BIC) son comúnmente utilizados.

21.5.1 Regresión Hacia Atrás

La regresión hacia atrás comienza con todas las covariables en el modelo y elimina iterativamente la covariable menos significativa hasta que todas las covariables restantes sean significativas.

21.5.2 Regresión Hacia Adelante

La regresión hacia adelante comienza con un modelo vacío y añade iterativamente la covariable más significativa hasta que no se pueda añadir ninguna covariable adicional significativa.

21.5.3 Criterios de Información

Los criterios de información, como el AIC (Akaike Information Criterion) y el BIC (Bayesian Information Criterion), se utilizan para seleccionar el modelo que mejor se ajusta a los datos con la menor complejidad posible:

$$\begin{aligned} AIC &= -2 \log L + 2k \\ BIC &= -2 \log L + k \log n \end{aligned}$$

donde L es la función de verosimilitud del modelo, k es el número de parámetros en el modelo y n es el tamaño de la muestra.

21.6 Ejemplo de Análisis Multivariado

Consideremos un ejemplo con tres covariables: edad, sexo y tratamiento. Ajustamos un modelo de Cox multivariado y obtenemos los siguientes coeficientes:

$$\hat{\beta}_{edad} = 0.03, \quad \hat{\beta}_{sexo} = -0.6, \quad \hat{\beta}_{tratamiento} = 1.5$$

La función de riesgo ajustada se expresa como:

$$\lambda(t | X) = \lambda_0(t) \exp(0.03 \cdot \text{edad} - 0.6 \cdot \text{sexo} + 1.5 \cdot \text{tratamiento})$$

21.7 Conclusión

El análisis multivariado de supervivencia permite evaluar el efecto conjunto de múltiples covariables sobre el tiempo hasta el evento. La inclusión de interacciones y efectos no lineales, junto con la selección adecuada de variables, mejora la precisión y la interpretabilidad de los modelos de supervivencia.

Chapter 22

Supervivencia en Datos Complicados

22.1 Introducción

El análisis de supervivencia en datos complicados se refiere a la evaluación de datos de supervivencia que presentan desafíos adicionales, como la censura por intervalo, datos truncados y datos con múltiples tipos de eventos. Estos escenarios requieren métodos avanzados para un análisis adecuado.

22.2 Censura por Intervalo

La censura por intervalo ocurre cuando el evento de interés se sabe que ocurrió dentro de un intervalo de tiempo, pero no se conoce el momento exacto. Esto es común en estudios donde las observaciones se realizan en puntos de tiempo discretos.

22.2.1 Modelo para Datos Censurados por Intervalo

Para datos censurados por intervalo, la función de verosimilitud se modifica para incluir la probabilidad de que el evento ocurra dentro de un intervalo:

$$L(\beta) = \prod_{i=1}^n P(T_i \in [L_i, U_i] \mid X_i; \beta)$$

donde $[L_i, U_i]$ es el intervalo de tiempo durante el cual se sabe que ocurrió el evento para el individuo i .

22.3 Datos Truncados

Los datos truncados ocurren cuando los tiempos de supervivencia están sujetos a un umbral, y solo se observan los individuos cuyos tiempos de supervivencia superan (o están por debajo de) ese umbral. Existen dos tipos principales de truncamiento: truncamiento a la izquierda y truncamiento a la derecha.

22.3.1 Modelo para Datos Truncados

Para datos truncados a la izquierda, la función de verosimilitud se ajusta para considerar solo los individuos que superan el umbral de truncamiento:

$$L(\beta) = \prod_{i=1}^n \frac{f(t_i \mid X_i; \beta)}{1 - F(L_i \mid X_i; \beta)}$$

donde L_i es el umbral de truncamiento para el individuo i .

22.4 Análisis de Competing Risks

En estudios donde pueden ocurrir múltiples tipos de eventos (competing risks), es crucial modelar adecuadamente el riesgo asociado con cada tipo de evento. La probabilidad de ocurrencia de cada evento compite con las probabilidades de ocurrencia de otros eventos.

22.4.1 Modelo de Competing Risks

Para un análisis de competing risks, la función de riesgo se descompone en funciones de riesgo específicas para cada tipo de evento:

$$\lambda(t) = \sum_{j=1}^m \lambda_j(t)$$

donde $\lambda_j(t)$ es la función de riesgo para el evento j .

22.5 Métodos de Imputación

Los métodos de imputación se utilizan para manejar datos faltantes o censurados en estudios de supervivencia. La imputación múltiple es un enfoque común que crea múltiples conjuntos de datos completos imputando valores faltantes varias veces y luego combina los resultados.

22.5.1 Imputación Múltiple

La imputación múltiple para datos de supervivencia se realiza en tres pasos:

1. Imputar los valores faltantes múltiples veces para crear varios conjuntos de datos completos.
2. Analizar cada conjunto de datos completo por separado utilizando métodos de supervivencia estándar.
3. Combinar los resultados de los análisis separados para obtener estimaciones y varianzas combinadas.

22.6 Ejemplo de Análisis con Datos Complicados

Consideremos un estudio con datos censurados por intervalo y competing risks. Ajustamos un modelo para los datos censurados por intervalo y obtenemos los siguientes coeficientes para las covariables edad y tratamiento:

$$\hat{\beta}_{edad} = 0.04, \quad \hat{\beta}_{tratamiento} = -0.8$$

La función de supervivencia ajustada se expresa como:

$$S(t | X) = \exp \left(- \left(\frac{t \exp(0.04 \cdot edad - 0.8 \cdot tratamiento)}{\lambda} \right)^k \right)$$

22.7 Conclusión

El análisis de supervivencia en datos complicados requiere métodos avanzados para manejar censura por intervalo, datos truncados y competing risks. La aplicación de modelos adecuados y métodos de imputación asegura un análisis preciso y completo de estos datos complejos.

Chapter 23

Proyecto Final y Revisión

23.1 Introducción

El proyecto final proporciona una oportunidad para aplicar los conceptos y técnicas aprendidas en el curso de análisis de supervivencia. Este capítulo incluye una guía para desarrollar un proyecto de análisis de supervivencia y una revisión de los conceptos clave.

23.2 Desarrollo del Proyecto

El proyecto final debe incluir los siguientes componentes:

1. Definición del problema: Identificar la pregunta de investigación y los objetivos del análisis de supervivencia.
2. Descripción de los datos: Presentar los datos utilizados, incluyendo las covariables y la estructura de los datos.
3. Análisis exploratorio: Realizar un análisis descriptivo de los datos, incluyendo la censura y la distribución de los tiempos de supervivencia.
4. Ajuste del modelo: Ajustar modelos de supervivencia adecuados (Kaplan-Meier, Cox, AFT) y evaluar su bondad de ajuste.
5. Diagnóstico del modelo: Realizar diagnósticos para evaluar los supuestos del modelo y la influencia de observaciones individuales.
6. Interpretación de resultados: Interpretar los coeficientes del modelo y las curvas de supervivencia ajustadas.
7. Conclusiones: Resumir los hallazgos del análisis y proporcionar recomendaciones basadas en los resultados.

23.3 Revisión de Conceptos Clave

Una revisión de los conceptos clave del análisis de supervivencia incluye:

- **Función de Supervivencia:** Define la probabilidad de sobrevivir más allá de un tiempo específico.
- **Función de Riesgo:** Define la tasa instantánea de ocurrencia del evento.

- **Estimador de Kaplan-Meier:** Proporciona una estimación no paramétrica de la función de supervivencia.
- **Test de Log-rank:** Compara curvas de supervivencia entre diferentes grupos.
- **Modelo de Cox:** Evalúa el efecto de múltiples covariables sobre el tiempo hasta el evento, asumiendo proporcionalidad de riesgos.
- **Modelos AFT:** Modelan el efecto de las covariables multiplicando el tiempo de supervivencia por una constante.
- **Análisis Multivariado:** Considera interacciones y efectos no lineales entre múltiples covariables.
- **Supervivencia en Datos Complicados:** Maneja censura por intervalo, datos truncados y competing risks.

23.4 Ejemplo de Proyecto Final

A continuación se presenta un ejemplo de estructura de un proyecto final de análisis de supervivencia:

23.4.1 Definición del Problema

Analizar el efecto del tratamiento y la edad sobre la supervivencia de pacientes con una enfermedad específica.

23.4.2 Descripción de los Datos

Datos de supervivencia de 100 pacientes, con covariables: edad, sexo y tipo de tratamiento. Los tiempos de supervivencia están censurados a la derecha.

23.4.3 Análisis Exploratorio

Realizar histogramas y curvas de Kaplan-Meier para explorar la distribución de los tiempos de supervivencia y la censura.

23.4.4 Ajuste del Modelo

Ajustar un modelo de Cox y un modelo AFT con las covariables edad y tratamiento.

23.4.5 Diagnóstico del Modelo

Evaluar la proporcionalidad de riesgos y realizar análisis de residuos para identificar observaciones influyentes.

23.4.6 Interpretación de Resultados

Interpretar los coeficientes del modelo y las curvas de supervivencia ajustadas para diferentes niveles de las covariables.

23.4.7 Conclusiones

Resumir los hallazgos y proporcionar recomendaciones sobre el efecto del tratamiento y la edad en la supervivencia de los pacientes.

23.5 Conclusión

El proyecto final es una oportunidad para aplicar los conocimientos adquiridos en un contexto práctico. La revisión de los conceptos clave y la aplicación de técnicas adecuadas de análisis de supervivencia aseguran un análisis riguroso y significativo.

Chapter 24

Fundamentos

24.1 2. Pruebas de Hipótesis

24.1.1 2.1 Tipos de errores

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado.
- El interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.
- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

Una prueba de hipótesis está formada por cinco partes:

- La hipótesis nula, denotada por H_0 .
- La hipótesis alternativa, denotada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 11. Las dos hipótesis en competencia son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones:

- Rechazar H_0 y concluir que H_1 es verdadera.
- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 5. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19 usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 6. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas:

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 12. El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadístico de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**.

Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**.

Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente.

Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 13. Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

24.2 2.2 Muestras grandes: una media poblacional

24.2.1 2.2.1 Cálculo de valor p

Definición 14. El **valor de p** (p -value) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.

Nota 18. Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.

Nota 19. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 15. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados:

p	H_0	Significativa
≤ 0.01	rechazada	Result. altamente significativos y en contra de H_0
≤ 0.05	rechazada	Result. significativos y en contra de H_0
≤ 0.10	rechazada	Result. posiblemente significativos y en contra de H_0
> 0.10	no rechazada	Result. no significativos y no rechazar H_0

Chapter 25

Elementos

25.1 Pruebas de Hipótesis

Tipos de errores

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado
- el interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.
- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

Una prueba de hipótesis está formada por cinco partes

- La hipótesis nula, denotada por H_0 .
- La hipótesis alterativa, denorada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 16. Las dos hipótesis en competencias son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones:

- Rechazar H_0 y concluir que H_1 es verdadera.
- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 7. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 8. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 17. El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadística de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**.

Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**.

Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente.

Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 18. Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

Muestras grandes: una media poblacional

Cálculo de valor p

Definición 19. El **valor de p** (p -value) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.

Nota 20. Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.

Nota 21. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 20. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados

p	H_0	Significativa
$p < 0.01$	Rechazar	Altamente
$0.01 \leq p < 0.05$	Rechazar	Estadísticamente
$0.05 \leq p < 0.1$	No rechazar	Tendencia estadística
$0.1 \leq p$	No rechazar	No son estadísticamente

Nota 22. Para determinar el valor de p , encontrar el área en la cola después del estadístico de prueba. Si la prueba es de una cola, este es el valor de p . Si es de dos colas, éste valor encontrado es la mitad del valor de p . Rechazar H_0 cuando el valor de $p < \alpha$.

Hay dos tipos de errores al realizar una prueba de hipótesis

	H_0 es Verdadera	H_0 es Falsa
Rechazar H_0	Error tipo I	✓
Aceptar H_0	✓	Error tipo II

Definición 21. La probabilidad de cometer el error tipo II se define por β donde

$$\begin{aligned}\beta &= P\{\text{error tipo II}\} = P\{\text{Aceptar equivocadamente } H_0\} \\ &= P\{\text{Aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\}\end{aligned}$$

Nota 23. Cuando H_0 es falsa y H_1 es verdadera, no siempre es posible especificar un valor exacto de μ , sino más bien un rango de posibles valores. En lugar de arriesgarse a tomar una decisión incorrecta, es mejor concluir que no hay evidencia suficiente para rechazar H_0 , es decir en lugar de aceptar H_0 , no rechazar H_0 .

La bondad de una prueba estadística se mide por el tamaño de α y β , ambas deben de ser pequeñas. Una manera muy efectiva de medir la potencia de la prueba es calculando el complemento del error tipo II:

$$\begin{aligned}1 - \beta &= P\{\text{Rechazar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{\text{Rechazar } H_0 \text{ cuando } H_1 \text{ es verdadera}\}\end{aligned}$$

Definición 22. La **potencia de la prueba**, $1 - \beta$, mide la capacidad de que la prueba funciona como se necesita.

Ejemplo 9. La producción diaria de una planta química local ha promediado 880 toneladas en los últimos años. A la gerente de control de calidad le gustaría saber si este promedio ha cambiado en meses recientes. Ella selecciona al azar 50 días de la base de datos computarizada y calcula el promedio y la desviación estándar de las $n = 50$ producciones como $\bar{x} = 871$ toneladas y $s = 21$ toneladas, respectivamente. Pruebe la hipótesis apropiada usando $\alpha = 0.05$.

La hipótesis nula apropiada es:

$$\begin{aligned}H_0 &: \mu = 880 \\ &\text{y la hipótesis alternativa } H_1 \text{ es} \\ H_1 &: \mu \neq 880\end{aligned}$$

el estimador puntual para μ es \bar{x} , entonces el estadístico de prueba es

$$\begin{aligned}
 z &= \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \\
 &= \frac{871 - 880}{21/\sqrt{50}} = -3.03
 \end{aligned}$$

Para esta prueba de dos colas, hay que determinar los dos valores de $z_{\alpha/2}$, es decir, $z_{\alpha/2} = \pm 1.96$, como $z > z_{\alpha/2}$, z cae en la zona de rechazo, por lo tanto la gerente puede rechazar la hipótesis nula y concluir que el promedio efectivamente ha cambiado. La probabilidad de rechazar H_0 cuando esta es verdadera es de 0.05.

Recordemos que el valor observado del estadístico de prueba es $z = -3.03$, la región de rechazo más pequeña que puede usarse y todavía seguir rechazando H_0 es $|z| > 3.03$, entonces $p = 2(0.012) = 0.0024$, que a su vez es menor que el nivel de significancia α asignado inicialmente, y además los resultados son **altamente significativos**.

Finalmente determinemos la potencia de la prueba cuando μ en realidad es igual a 870 toneladas.

Recordar que la región de aceptación está entre -1.96 y 1.96 , para $\mu = 880$, equivalentemente

$$874.18 < \bar{x} < 885.82$$

β es la probabilidad de aceptar H_0 cuando $\mu = 870$, calculemos los valores de z correspondientes a 874.18 y 885.82. Entonces

$$\begin{aligned}
 z_1 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{874.18 - 870}{21/\sqrt{50}} = 1.41 \\
 z_2 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{885.82 - 870}{21/\sqrt{50}} = 5.33
 \end{aligned}$$

por lo tanto

$$\begin{aligned}
 \beta &= P\{\text{aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\
 &= P\{874.18 < \mu < 885.82 \text{ cuando } \mu = 870\} \\
 &= P\{1.41 < z < 5.33\} = P\{1.41 < z\} \\
 &= 1 - 0.9207 = 0.0793
 \end{aligned}$$

entonces, la potencia de la prueba es

$$1 - \beta = 1 - 0.0793 = 0.9207$$

que es la probabilidad de rechazar correctamente H_0 cuando H_0 es falsa.

Prueba de hipótesis para la diferencia entre dos medias poblacionales

El estadístico que resume la información muestral respecto a la diferencia en medias poblacionales ($\mu_1 - \mu_2$) es la diferencia de las medias muestrales ($\bar{x}_1 - \bar{x}_2$), por tanto al probar la diferencia entre las medias muestrales se verifica que la diferencia real entre las medias poblacionales difiere de un valor especificado, $(\mu_1 - \mu_2) = D_0$, se puede usar el error estándar de ($\bar{x}_1 - \bar{x}_2$), es decir

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

cuyo estimador está dado por

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

El procedimiento para muestras grandes es:

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$, donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\mu_1 - \mu_2) > D_0$ $H_1 : (\mu_1 - \mu_2) < D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$

- 3) Estadístico de prueba:

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_0$ $z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Ejemplo 10. Para determinar si ser propietario de un automóvil afecta el rendimiento académico de un estudiante, se tomaron dos muestras aleatorias de 100 estudiantes varones. El promedio de calificaciones para los $n_1 = 100$ no propietarios de un auto tuvieron un promedio y varianza de $\bar{x}_1 = 2.7$ y $s_1^2 = 0.36$, respectivamente, mientras que para la segunda muestra con $n_2 = 100$ propietarios de un auto, se tiene $\bar{x}_2 = 2.54$ y $s_2^2 = 0.4$. Los datos presentan suficiente evidencia para indicar una diferencia en la media en el rendimiento académico entre propietarios y no propietarios de un automóvil? Hacer pruebas para $\alpha = 0.01, 0.05$ y $\alpha = 0.1$.

- Solución utilizando la técnica de regiones de rechazo: realizando las operaciones $z = 1.84$, determinar si excede los valores de $z_{\alpha/2}$.
- Solución utilizando el p-value: Calcular el valor de p , la probabilidad de que z sea mayor que $z = 1.84$ o menor que $z = -1.84$, se tiene que $p = 0.0658$. Concluir.
- Si el intervalo de confianza que se construye contiene el valor del parámetro especificado por H_0 , entonces ese valor es uno de los posibles valores del parámetro y H_0 no debe ser rechazada.
- Si el valor hipotético se encuentra fuera de los límites de confianza, la hipótesis nula es rechazada al nivel de significancia α .

Prueba de Hipótesis para una Proporción Binomial

Para una muestra aleatoria de n intentos idénticos, de una población binomial, la proporción muestral \hat{p} tiene una distribución aproximadamente normal cuando n es grande, con media p y error estándar

$$SE = \sqrt{\frac{pq}{n}}.$$

La prueba de hipótesis de la forma

$$\begin{aligned} H_0 &: p = p_0 \\ H_1 &: p > p_0, \text{ o } p < p_0 \text{ o } p \neq p_0 \end{aligned}$$

El estadístico de prueba se construye con el mejor estimador de la proporción verdadera, \hat{p} , con el estadístico de prueba z , que se distribuye normal estándar.

El procedimiento es

- 1) Hipótesis nula: $H_0 : p = p_0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis alternativa	$H_1 : p > p_0$ $H_1 : p < p_0$	$p \neq p_0$

3) Estadístico de prueba:

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{pq}{n}}}, \hat{p} = \frac{x}{n}$$

donde x es el número de éxitos en n intentos binomiales.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_0$ $z < -z_\alpha$ cuando $H_1 : p < p_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Prueba de Hipótesis diferencia entre dos Proporciones Binomiales

Cuando se tienen dos muestras aleatorias independientes de dos poblaciones binomiales, el objetivo del experimento puede ser la diferencia $(p_1 - p_2)$ en las proporciones de individuos u objetos que poseen una característica específica en las dos poblaciones. En este caso se pueden utilizar los estimadores de las dos proporciones $(\hat{p}_1 - \hat{p}_2)$ con error estándar dado por

$$SE = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

considerando el estadístico z con un nivel de significancia $(1 - \alpha) 100\%$

La hipótesis nula a probarse es de la forma

H_0 : $p_1 = p_2$ o equivalentemente $(p_1 - p_2) = 0$, contra una hipótesis alternativa H_1 de una o dos colas.

Para estimar el error estándar del estadístico z , se debe de utilizar el hecho de que suponiendo que H_0 es verdadera, las dos proporciones son iguales a algún valor común, p . Para obtener el mejor estimador de p es

$$p = \frac{\text{número total de éxitos}}{\text{Número total de pruebas}} = \frac{x_1 + x_2}{n_1 + n_2}$$

1) **Hipótesis Nula:** $H_0 : (p_1 - p_2) = 0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: H_1 :	$H_1 : (p_1 - p_2) > 0$ $H_1 : (p_1 - p_2) < 0$	$H_1 : (p_1 - p_2) \neq 0$

3) Estadístico de prueba:

$$z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{pq}{n_1} + \frac{pq}{n_2}}}$$

donde $\hat{p}_1 = x_1/n_1$ y $\hat{p}_2 = x_2/n_2$, dado que el valor común para p_1 y p_2 es p , entonces $\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$ y por tanto el estadístico de prueba es

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_\alpha$ $z < -z_\alpha$ cuando $H_1 : p < p_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

25.2 Muestras Pequeñas

Una media poblacional

- 1) **Hipótesis Nula:** $H_0 : \mu = \mu_0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 :$	$H_1 : \mu > \mu_0$ $H_1 : \mu < \mu_0$	$H_1 : \mu \neq \mu_0$

- 3) Estadístico de prueba:

$$t = \frac{\bar{x} - \mu_0}{\sqrt{\frac{s^2}{n}}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Diferencia entre dos medias poblacionales: Muestras Aleatorias Independientes

Cuando los tamaños de muestra son pequeños, no se puede asegurar que las medias muestrales sean normales, pero si las poblaciones originales son normales, entonces la distribución muestral de la diferencia de las medias muestrales, $(\bar{x}_1 - \bar{x}_2)$, será normal con media $(\mu_1 - \mu_2)$ y error estándar

$$ES = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\mu_1 - \mu_2) > D_0$ $H_1 : (\mu_1 - \mu_2) < D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$

- 3) Estadístico de prueba:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

donde

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

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_\alpha$ $z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

Diferencia entre dos medias poblacionales: Diferencias Pareadas

- 1) **Hipótesis Nula:** $H_0 : \mu_d = 0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 : \mu_d$	$H_1 : \mu_d > 0$ $H_1 : \mu_d < 0$	$H_1 : \mu_d \neq 0$

- 3) Estadístico de prueba:

$$t = \frac{\bar{d}}{\sqrt{\frac{s_d^2}{n}}}$$

donde n es el número de diferencias pareadas, \bar{d} es la media de las diferencias muestrales, y s_d es la desviación estándar de las diferencias muestrales.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

Inferencias con respecto a la Varianza Poblacional

- 1) **Hipótesis Nula:** $H_0 : \sigma^2 = \sigma_0^2$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: H_1	$H_1 : \sigma^2 > \sigma_0^2$ $H_1 : \sigma^2 < \sigma_0^2$	$H_1 : \sigma^2 \neq \sigma_0^2$

- 3) Estadístico de prueba:

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$\chi^2 > \chi_\alpha^2$ $\chi^2 < \chi_{(1-\alpha)}^2$ cuando $H_1 : \chi^2 < \chi_0^2$ cuando $p < \alpha$	$\chi^2 > \chi_{\alpha/2}^2$ o $\chi^2 < \chi_{(1-\alpha/2)}^2$

Los valores críticos de χ^2 , están basados en $(n_1 +)$ grados de libertad.

Comparación de dos varianzas poblacionales

- 1) **Hipótesis Nula** $H_0 : (\sigma_1^2 - \sigma_2^2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\sigma_1^2 - \sigma_2^2) > D_0$ $H_1 : (\sigma_1^2 - \sigma_2^2) < D_0$	$H_1 : (\sigma_1^2 - \sigma_2^2) \neq D_0$

3) Estadístico de prueba:

$$F = \frac{s_1^2}{s_2^2}$$

donde s_1^2 es la varianza muestral más grande.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$F > F_\alpha$ cuando $p < \alpha$	$F > F_{\alpha/2}$

25.3 Introducción

La regresión logística es una técnica de modelado estadístico utilizada para predecir la probabilidad de un evento binario en función de una o más variables independientes. Este capítulo profundiza en las matemáticas subyacentes a la regresión logística, incluyendo la función logística, la función de verosimilitud, y los métodos para estimar los coeficientes del modelo.

25.4 Función Logística

La función logística es la base de la regresión logística. Esta función transforma una combinación lineal de variables independientes en una probabilidad.

25.4.1 Definición

La función logística se define como:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$

donde p es la probabilidad de que el evento ocurra, $\beta_0, \beta_1, \dots, \beta_n$ son los coeficientes del modelo, y X_1, X_2, \dots, X_n son las variables independientes.

25.4.2 Propiedades

La función logística tiene varias propiedades importantes:

- **Rango:** La función logística siempre produce un valor entre 0 y 1, lo que la hace adecuada para modelar probabilidades.
- **Monotonía:** La función es monótona creciente, lo que significa que a medida que la combinación lineal de variables independientes aumenta, la probabilidad también aumenta.
- **Simetría:** La función logística es simétrica en torno a $p = 0.5$.

25.5 Función de Verosimilitud

La función de verosimilitud se utiliza para estimar los coeficientes del modelo de regresión logística. Esta función mide la probabilidad de observar los datos dados los coeficientes del modelo.

25.5.1 Definición

Para un conjunto de n observaciones, la función de verosimilitud L se define como el producto de las probabilidades individuales de observar cada dato:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación y p_i es la probabilidad predicha de que $Y_i = 1$.

25.5.2 Función de Log-Verosimilitud

Para simplificar los cálculos, trabajamos con el logaritmo de la función de verosimilitud, conocido como la función de log-verosimilitud. Tomar el logaritmo convierte el producto en una suma:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

Sustituyendo p_i :

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}) \right]$$

25.6 Estimación de Coeficientes

Los coeficientes del modelo de regresión logística se estiman maximizando la función de log-verosimilitud. Este proceso generalmente se realiza mediante métodos iterativos como el algoritmo de Newton-Raphson.

25.6.1 Gradiente y Hessiana

Para maximizar la función de log-verosimilitud, necesitamos calcular su gradiente y su matriz Hessiana.

Gradiente

El gradiente de la función de log-verosimilitud con respecto a los coeficientes β es:

$$\mathbf{g}(\beta) = \frac{\partial \log L}{\partial \beta} = \sum_{i=1}^n \mathbf{X}_i (y_i - p_i)$$

donde \mathbf{X}_i es el vector de valores de las variables independientes para la i -ésima observación.

Hessiana

La matriz Hessiana de la función de log-verosimilitud con respecto a los coeficientes β es:

$$\mathbf{H}(\beta) = \frac{\partial^2 \log L}{\partial \beta \partial \beta^T} = - \sum_{i=1}^n p_i (1 - p_i) \mathbf{X}_i \mathbf{X}_i^T$$

25.6.2 Algoritmo Newton-Raphson

El algoritmo Newton-Raphson se utiliza para encontrar los valores de los coeficientes que maximizan la función de log-verosimilitud. El algoritmo se puede resumir en los siguientes pasos:

1. Inicializar el vector de coeficientes $\beta^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).
2. Calcular el gradiente $\mathbf{g}(\beta^{(k)})$ y la matriz Hessiana $\mathbf{H}(\beta^{(k)})$ en la iteración k .
3. Actualizar los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - [\mathbf{H}(\beta^{(k)})]^{-1} \mathbf{g}(\beta^{(k)})$$

4. Repetir los pasos 2 y 3 hasta que la diferencia entre $\beta^{(k+1)}$ y $\beta^{(k)}$ sea menor que un umbral predefinido (criterio de convergencia).

25.7 Validación del Modelo

Una vez que se han estimado los coeficientes del modelo de regresión logística, es importante validar el modelo para asegurarse de que proporciona predicciones precisas.

25.7.1 Curva ROC y AUC

La curva ROC (Receiver Operating Characteristic) es una herramienta gráfica utilizada para evaluar el rendimiento de un modelo de clasificación binaria. El área bajo la curva (AUC) mide la capacidad del modelo para distinguir entre las clases.

25.7.2 Matriz de Confusión

La matriz de confusión es una tabla que resume el rendimiento de un modelo de clasificación al comparar las predicciones del modelo con los valores reales. Los términos en la matriz de confusión incluyen verdaderos positivos, falsos positivos, verdaderos negativos y falsos negativos.

25.8 Conceptos Básicos

La regresión logística es una técnica de modelado estadístico utilizada para predecir la probabilidad de un evento binario (es decir, un evento que tiene dos posibles resultados) en función de una o más variables independientes. Es ampliamente utilizada en diversas disciplinas, como medicina, economía, biología, y ciencias sociales, para analizar y predecir resultados binarios. Un modelo de regresión logística describe cómo una variable dependiente binaria Y (que puede tomar los valores 0 o 1) está relacionada con una o más variables independientes X_1, X_2, \dots, X_n . A diferencia de la regresión lineal, que predice un valor continuo, la regresión logística predice una probabilidad que puede ser interpretada como la probabilidad de que $Y = 1$ dado un conjunto de valores para X_1, X_2, \dots, X_n .

25.9 Regresión Lineal

La regresión lineal es utilizada para predecir el valor de una variable dependiente continua en función de una o más variables independientes. El modelo de regresión lineal tiene la forma:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon \quad (25.1)$$

donde:

- Y es la variable dependiente.
- β_0 es la intersección con el eje Y o término constante.
- $\beta_1, \beta_2, \dots, \beta_n$ son los coeficientes que representan la relación entre las variables independientes y la variable dependiente.
- X_1, X_2, \dots, X_n son las variables independientes.
- ϵ es el término de error, que representa la desviación de los datos observados de los valores predichos por el modelo.

El objetivo de la regresión lineal es encontrar los valores de los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ que minimicen la suma de los cuadrados de las diferencias entre los valores observados y los valores predichos. Este método se conoce como mínimos cuadrados ordinarios (OLS, por sus siglas en inglés). La función de costo a minimizar es:

$$J(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (25.2)$$

donde:

- y_i es el valor observado de la variable dependiente para la i -ésima observación.
- \hat{y}_i es el valor predicho por el modelo para la i -ésima observación, dado por:

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_n x_{in} \quad (25.3)$$

Para encontrar los valores óptimos de los coeficientes, se toman las derivadas parciales de la función de costo con respecto a cada coeficiente y se igualan a cero:

$$\frac{\partial J}{\partial \beta_j} = 0 \quad \text{para } j = 0, 1, \dots, n \quad (25.4)$$

Resolviendo este sistema de ecuaciones, se obtienen los valores de los coeficientes que minimizan la función de costo.

25.10 Regresión Logística

La deducción de la fórmula de la regresión logística comienza con la necesidad de modelar la probabilidad de un evento binario. Queremos encontrar una función que relacione las variables independientes con la probabilidad de que la variable dependiente tome el valor 1. La probabilidad de que el evento ocurra, $P(Y = 1)$, se denota como p . La probabilidad de que el evento no ocurra, $P(Y = 0)$, es $1 - p$. Los *odds* (chances) de que ocurra el evento se definen como:

$$\text{odds} = \frac{p}{1 - p} \quad (25.5)$$

Los *odds* nos indican cuántas veces más probable es que ocurra el evento frente a que no ocurra. Para simplificar el modelado de los *odds*, aplicamos el logaritmo natural, obteniendo la función logit:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) \quad (25.6)$$

La transformación logit es útil porque convierte el rango de la probabilidad $(0, 1)$ al rango de números reales $(-\infty, \infty)$. La idea clave de la regresión logística es modelar la transformación logit de la probabilidad como una combinación lineal de las variables independientes:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad (25.7)$$

Aquí, β_0 es el término constante y $\beta_1, \beta_2, \dots, \beta_n$ son los coeficientes asociados con las variables independientes X_1, X_2, \dots, X_n . Para expresar p en función de una combinación lineal de las variables independientes, invertimos la transformación logit. Partimos de la ecuación:

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

Aplicamos la exponenciación a ambos lados:

$$\frac{p}{1-p} = e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n}$$

Despejamos p :

$$p = \frac{e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n}}{1 + e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n}}$$

La expresión final que obtenemos es conocida como la función logística:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} \quad (25.8)$$

Esta función describe cómo las variables independientes se relacionan con la probabilidad de que el evento de interés ocurra. Los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ se estiman a partir de los datos utilizando el método de máxima verosimilitud.

25.11 Método de Máxima Verosimilitud

Para estimar los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ en la regresión logística, utilizamos el método de máxima verosimilitud. La idea es encontrar los valores de los coeficientes que maximicen la probabilidad de observar los datos dados. Esta probabilidad se expresa mediante la función de verosimilitud L . La función de verosimilitud $L(\beta_0, \beta_1, \dots, \beta_n)$ para un conjunto de n observaciones se define como el producto de las probabilidades de las observaciones dadas las variables independientes:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (25.9)$$

donde:

- p_i es la probabilidad predicha de que $Y_i = 1$,
- y_i es el valor observado de la variable dependiente para la i -ésima observación.

Trabajar directamente con esta función de verosimilitud puede ser complicado debido al producto de muchas probabilidades, especialmente si n es grande. Para simplificar los cálculos, se utiliza el logaritmo de la función de verosimilitud, conocido como la función de log-verosimilitud. El uso del logaritmo simplifica significativamente la diferenciación y maximización de la función. La función de log-verosimilitud se define como:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (25.10)$$

Aquí, \log representa el logaritmo natural. Esta transformación es válida porque el logaritmo es una función monótona creciente, lo que significa que maximizar la log-verosimilitud es equivalente a maximizar la verosimilitud original. En la regresión logística, la probabilidad p_i está dada por la función logística:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \quad (25.11)$$

Sustituyendo esta expresión en la función de log-verosimilitud, obtenemos:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n \left[y_i \log \left(\frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \right) + \right. \\ &\quad \left. (1 - y_i) \log \left(1 - \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \right) \right] \end{aligned}$$

Simplificando esta expresión, notamos que:

$$\log \left(\frac{1}{1 + e^{-z}} \right) = -\log(1 + e^{-z})$$

y

$$\log \left(1 - \frac{1}{1 + e^{-z}} \right) = \log \left(\frac{e^{-z}}{1 + e^{-z}} \right) = -z - \log(1 + e^{-z})$$

Aplicando estas identidades, la función de log-verosimilitud se convierte en:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n \left[y_i (-\log(1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})})) + \right. \\ &\quad \left. (1 - y_i) \left(-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}) \right) \right] \end{aligned}$$

Simplificando aún más, obtenemos:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) \right. \\ &\quad \left. - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}) \right] \end{aligned}$$

Para simplificar aún más la notación, podemos utilizar notación matricial. Definimos la matriz \mathbf{X} de tamaño $n \times (k + 1)$ y el vector de coeficientes $\boldsymbol{\beta}$ de tamaño $(k + 1) \times 1$ como sigue:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \quad (25.12)$$

Entonces, la expresión para la función de log-verosimilitud es:

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n \left[y_i(\mathbf{X}_i \boldsymbol{\beta}) - \log(1 + e^{\mathbf{X}_i \boldsymbol{\beta}}) \right] \quad (25.13)$$

donde \mathbf{X}_i es la i -ésima fila de la matriz \mathbf{X} . Esta notación matricial simplifica la implementación y la derivación de los estimadores de los coeficientes en la regresión logística. Utilizando métodos numéricos, como el algoritmo de Newton-Raphson, se pueden encontrar los coeficientes que maximizan la función de log-verosimilitud. Para maximizar la función de log-verosimilitud, derivamos esta función con respecto a cada uno de los coeficientes β_j y encontramos los puntos críticos. La derivada parcial de la función de log-verosimilitud con respecto a β_j es:

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = \sum_{i=1}^n \left[y_i X_{ij} - \frac{X_{ij} e^{\mathbf{X}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{X}_i \boldsymbol{\beta}}} \right] \quad (25.14)$$

Simplificando, esta derivada se puede expresar como:

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = \sum_{i=1}^n X_{ij}(y_i - p_i), \quad \text{donde } p_i = \frac{1}{1 + e^{-\mathbf{X}_i \boldsymbol{\beta}}} \quad (25.15)$$

Para encontrar los coeficientes que maximizan la log-verosimilitud, resolvemos el sistema de ecuaciones

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = 0 \quad \text{para todos los } j = 0, 1, \dots, k.$$

Este sistema de ecuaciones no tiene una solución analítica cerrada, por lo que se resuelve numéricamente utilizando métodos iterativos como el algoritmo de Newton-Raphson.

25.12 Método de Newton-Raphson

El método de Newton-Raphson es un algoritmo iterativo que se utiliza para encontrar las raíces de una función. En el contexto de la regresión logística, se utiliza para maximizar la función de log-verosimilitud encontrando los valores de los coeficientes $\beta_0, \beta_1, \dots, \beta_n$. Este método se basa en una aproximación de segundo orden de la función objetivo. Dado un valor inicial de los coeficientes $\boldsymbol{\beta}^{(0)}$, se actualiza iterativamente el valor de los coeficientes utilizando la fórmula:

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \left[\mathbf{H}(\boldsymbol{\beta}^{(t)}) \right]^{-1} \nabla \log L(\boldsymbol{\beta}^{(t)}) \quad (25.16)$$

donde:

- $\beta^{(t)}$ es el vector de coeficientes en la t -ésima iteración.
- $\nabla \log L(\beta^{(t)})$ es el gradiente de la función de log-verosimilitud con respecto a los coeficientes β :

$$\nabla \log L(\beta) = \mathbf{X}^T(\mathbf{y} - \mathbf{p}) \quad (25.17)$$

donde \mathbf{y} es el vector de valores observados y \mathbf{p} es el vector de probabilidades.

- $\mathbf{H}(\beta^{(t)})$ es la matriz Hessiana (matriz de segundas derivadas) evaluada en $\beta^{(t)}$:

$$\mathbf{H}(\beta) = -\mathbf{X}^T \mathbf{W} \mathbf{X} \quad (25.18)$$

donde \mathbf{W} es una matriz diagonal de pesos con elementos $w_i = p_i(1 - p_i)$.

En resumen:

Algoritmo 3. *El algoritmo Newton-Raphson para la regresión logística se puede resumir en los siguientes pasos:*

1. Inicializar el vector de coeficientes $\beta^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).
2. Calcular el gradiente $\nabla \log L(\beta^{(t)})$ y la matriz Hessiana $\mathbf{H}(\beta^{(t)})$ en la iteración t .
3. Actualizar los coeficientes utilizando la fórmula:

$$\beta^{(t+1)} = \beta^{(t)} - [\mathbf{H}(\beta^{(t)})]^{-1} \nabla \log L(\beta^{(t)}) \quad (25.19)$$

4. Repetir los pasos 2 y 3 hasta que la diferencia entre $\beta^{(t+1)}$ y $\beta^{(t)}$ sea menor que un umbral pre-definido (criterio de convergencia).

En resumen, el método de Newton-Raphson permite encontrar los coeficientes que maximizan la función de log-verosimilitud de manera eficiente.

25.13 Especificando

En específico para un conjunto de n observaciones, la función de verosimilitud L se define como el producto de las probabilidades individuales de observar cada dato:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (25.20)$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación y p_i es la probabilidad predicha de que $Y_i = 1$. Aquí, p_i es dado por la función logística:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \quad (25.21)$$

Tomando el logaritmo:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (25.22)$$

Sustituyendo p_i :

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}) \right] \quad (25.23)$$

Dado que el objetivo es encontrar los valores de $\beta_0, \beta_1, \dots, \beta_n$ que maximicen la función de log-verosimilitud. Para β_j , la derivada parcial de la función de log-verosimilitud es:

$$\frac{\partial \log L}{\partial \beta_j} = \sum_{i=1}^n \left[y_i X_{ij} - \frac{X_{ij} e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}}{1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}} \right] \quad (25.24)$$

Esto se simplifica a (comparar con la ecuación ??):

$$\frac{\partial \log L}{\partial \beta_j} = \sum_{i=1}^n X_{ij} (y_i - p_i) \quad (25.25)$$

Para maximizar la log-verosimilitud, resolvemos el sistema de ecuaciones $\frac{\partial \log L}{\partial \beta_j} = 0$ para todos los j de 0 a n , mismo que se resuelve numéricamente utilizando métodos el algoritmo de Newton-Raphson. El método de Newton-Raphson se basa en una aproximación de segundo orden de la función objetivo. Dado un valor inicial de los coeficientes $\beta^{(0)}$, se iterativamente actualiza el valor de los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - [\mathbf{H}(\beta^{(k)})]^{-1} \mathbf{g}(\beta^{(k)}) \quad (25.26)$$

donde:

- $\beta^{(k)}$ es el vector de coeficientes en la k -ésima iteración.
- $\mathbf{g}(\beta^{(k)})$ es el gradiente (vector de primeras derivadas) evaluado en $\beta^{(k)}$:

$$\mathbf{g}(\beta) = \frac{\partial \log L}{\partial \beta} = \sum_{i=1}^n \mathbf{X}_i (y_i - p_i) \quad (25.27)$$

donde \mathbf{X}_i es el vector de valores de las variables independientes para la i -ésima observación (comparar con ecuación ??).

- $\mathbf{H}(\beta^{(k)})$ es la matriz Hessiana (matriz de segundas derivadas) evaluada en $\beta^{(k)}$:

$$\mathbf{H}(\beta) = \frac{\partial^2 \log L}{\partial \beta \partial \beta^T} = - \sum_{i=1}^n p_i (1 - p_i) \mathbf{X}_i \mathbf{X}_i^T, \quad (25.28)$$

comparar con ecuación ??

Algoritmo 4. Los pasos del algoritmo Newton-Raphson para la regresión logística son:

1. Inicializar el vector de coeficientes $\beta^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).
2. Calcular el gradiente $\mathbf{g}(\beta^{(k)})$ y la matriz Hessiana $\mathbf{H}(\beta^{(k)})$ en la iteración k .
3. Actualizar los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - [\mathbf{H}(\beta^{(k)})]^{-1} \mathbf{g}(\beta^{(k)}) \quad (25.29)$$

4. Repetir los pasos 2 y 3 hasta que la diferencia entre $\beta^{(k+1)}$ y $\beta^{(k)}$ sea menor que un umbral pre-definido (criterio de convergencia).

Como se puede observar la diferencia entre el Algoritmo ?? y el Algoritmo ?? son mínimas

Notas finales

En el contexto de la regresión logística, los vectores X_1, X_2, \dots, X_n representan las variables independientes. Cada X_j es un vector columna que contiene los valores de la variable independiente j para cada una de las n observaciones. Es decir,

$$X_j = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{bmatrix} \quad (25.30)$$

Para simplificar la notación y los cálculos, a menudo combinamos todos los vectores de variables independientes en una única matriz de diseño \mathbf{X} de tamaño $n \times (k + 1)$, donde n es el número de observaciones y $k + 1$ es el número de variables independientes más el término de intercepto. La primera columna de \mathbf{X} corresponde a un vector de unos para el término de intercepto, y las demás columnas corresponden a los valores de las variables independientes:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix} \quad (25.31)$$

revisar la ecuación ???. De esta forma, el modelo logit puede ser escrito de manera compacta utilizando la notación matricial:

$$\text{logit}(p) = \log \left(\frac{p}{1 - p} \right) = \mathbf{X}\boldsymbol{\beta} \quad (25.32)$$

donde $\boldsymbol{\beta}$ es el vector de coeficientes:

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \quad (25.33)$$

Así, la probabilidad p se puede expresar como:

$$p = \frac{1}{1 + e^{-\mathbf{X}\boldsymbol{\beta}}} \quad (25.34)$$

Comparar la ecuación anterior con la ecuación ??. Esta notación matricial simplifica la implementación y la derivación de los estimadores de los coeficientes en la regresión logística. Para estimar los coeficientes $\boldsymbol{\beta}$ en la regresión logística, se utiliza el método de máxima verosimilitud. La función de verosimilitud $L(\boldsymbol{\beta})$ se define como el producto de las probabilidades de las observaciones dadas las variables independientes, recordemos la ecuación ??:

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (25.35)$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación, y p_i es la probabilidad predicha de que $Y_i = 1$. La función de log-verosimilitud, que es más fácil de maximizar, se obtiene tomando el logaritmo natural de la función de verosimilitud (??):

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (25.36)$$

Sustituyendo $p_i = \frac{1}{1+e^{-\mathbf{X}_i\boldsymbol{\beta}}}$, donde \mathbf{X}_i es la i -ésima fila de la matriz de diseño \mathbf{X} , obtenemos:

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n \left[y_i(\mathbf{X}_i\boldsymbol{\beta}) - \log(1 + e^{\mathbf{X}_i\boldsymbol{\beta}}) \right] \quad (25.37)$$

Para encontrar los valores de $\boldsymbol{\beta}$ que maximizan la función de log-verosimilitud, se utiliza un algoritmo iterativo como el método de Newton-Raphson. Este método requiere calcular el gradiente y la matriz Hessiana de la función de log-verosimilitud.

El gradiente de la función de log-verosimilitud con respecto a $\boldsymbol{\beta}$ es (?? y ??):

$$\nabla \log L(\boldsymbol{\beta}) = \mathbf{X}^T(\mathbf{y} - \mathbf{p}) \quad (25.38)$$

donde \mathbf{y} es el vector de valores observados y \mathbf{p} es el vector de probabilidades predichas.

La matriz Hessiana de la función de log-verosimilitud es (?? y ??):

$$\mathbf{H}(\boldsymbol{\beta}) = -\mathbf{X}^T \mathbf{W} \mathbf{X} \quad (25.39)$$

donde \mathbf{W} es una matriz diagonal de pesos con elementos $w_i = p_i(1 - p_i)$.

El método de Newton-Raphson actualiza los coeficientes $\boldsymbol{\beta}$ de la siguiente manera:

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - [\mathbf{H}(\boldsymbol{\beta}^{(t)})]^{-1} \nabla \log L(\boldsymbol{\beta}^{(t)}) \quad (25.40)$$

Iterando este proceso hasta que la diferencia entre $\boldsymbol{\beta}^{(t+1)}$ y $\boldsymbol{\beta}^{(t)}$ sea menor que un umbral predefinido (??, ??, ?? y ??), se obtienen los estimadores de máxima verosimilitud para los coeficientes de la regresión logística.

Chapter 26

Introducción

26.1 Introducción

La Estadística es una ciencia formal que estudia la recolección, análisis e interpretación de datos de una muestra representativa, ya sea para ayudar en la toma de decisiones o para explicar condiciones regulares o irregulares de algún fenómeno o estudio aplicado, de ocurrencia en forma aleatoria o condicional. Sin embargo, la estadística es más que eso, es decir, es transversal a una amplia variedad de disciplinas, desde la física hasta las ciencias sociales, desde las ciencias de la salud hasta el control de calidad. Se usa para la toma de decisiones en áreas de negocios o instituciones gubernamentales. Ahora bien, las técnicas estadísticas se aplican de manera amplia en mercadotecnia, contabilidad, control de calidad y en otras actividades; estudios de consumidores; análisis de resultados en deportes; administradores de instituciones; en la educación; organismos políticos; médicos; y por otras personas que intervienen en la toma de decisiones.

Definición 23. *La Estadística es la ciencia cuyo objetivo es reunir una información cuantitativa concerniente a individuos, grupos, series de hechos, etc. y deducir de ello gracias al análisis de estos datos unos significados precisos o unas previsiones para el futuro.*

La estadística, en general, es la ciencia que trata de la recopilación, organización presentación, análisis e interpretación de datos numéricos con el fin de realizar una toma de decisión más efectiva. Los métodos estadísticos tradicionalmente se utilizan para propósitos descriptivos, para organizar y resumir datos numéricos. La estadística descriptiva, por ejemplo trata de la tabulación de datos, su presentación en forma gráfica o ilustrativa y el cálculo de medidas descriptivas.

26.1.1 Historia de la Estadística

Es difícil conocer los orígenes de la Estadística. Desde los comienzos de la civilización han existido formas sencillas de estadística, pues ya se utilizaban representaciones gráficas y otros símbolos en pieles, rocas, palos de madera y paredes de cuevas para contar el número de personas, animales o ciertas cosas. Su origen empieza posiblemente en la isla de Cerdeña, donde existen monumentos prehistóricos pertenecientes a los Nuragas, los primeros habitantes de la isla; estos monumentos constan de bloques de basalto superpuestos sin mortero y en cuyas paredes se encontraban grabados toscos signos que han sido interpretados con mucha verosimilitud como muescas que servían para llevar la cuenta del ganado y la caza. Los babilonios usaban ya pequeñas tablillas de arcilla para recopilar datos en tablas sobre la producción agrícola y los géneros vendidos o cambiados mediante trueque. Otros vestigios pueden ser hallados en el antiguo Egipto, cuyos faraones lograron recopilar, hacia el año 3050 antes de Cristo, prolijos datos relativos a la población y la riqueza del país. De acuerdo al historiador griego Heródoto,

dicho registro de riqueza y población se hizo con el objetivo de preparar la construcción de las pirámides. En el mismo Egipto, Ramsés II hizo un censo de las tierras con el objeto de verificar un nuevo reparto. En el antiguo Israel la Biblia da referencias, en el libro de los Números, de los datos estadísticos obtenidos en dos recuentos de la población hebrea. El rey David por otra parte, ordenó a Joab, general del ejército hacer un censo de Israel con la finalidad de conocer el número de la población. También los chinos efectuaron censos hace más de cuarenta siglos. Los griegos efectuaron censos periódicamente con fines tributarios, sociales (división de tierras) y militares (cálculo de recursos y hombres disponibles). La investigación histórica revela que se realizaron 69 censos para calcular los impuestos, determinar los derechos de voto y ponderar la potencia guerrera.

Fueron los romanos, maestros de la organización política, quienes mejor supieron emplear los recursos de la estadística. Cada cinco años realizaban un censo de la población y sus funcionarios públicos tenían la obligación de anotar nacimientos, defunciones y matrimonios, sin olvidar los recuentos periódicos del ganado y de las riquezas contenidas en las tierras conquistadas. Para el nacimiento de Cristo sucedía uno de estos empadronamientos de la población bajo la autoridad del imperio. Durante los mil años siguientes a la caída del imperio Romano se realizaron muy pocas operaciones Estadísticas, con la notable excepción de las relaciones de tierras pertenecientes a la Iglesia, compiladas por Pipino el Breve en el 758 y por Carlomagno en el 762 DC. Durante el siglo IX se realizaron en Francia algunos censos parciales de siervos. En Inglaterra, Guillermo el Conquistador recopiló el Domesday Book o libro del Gran Catastro para el año 1086, un documento de la propiedad, extensión y valor de las tierras de Inglaterra. Esa obra fue el primer compendio estadístico de Inglaterra. Aunque Carlomagno, en Francia; y Guillermo el Conquistador, en Inglaterra, trataron de revivir la técnica romana, los métodos estadísticos permanecieron casi olvidados durante la Edad Media.

Durante los siglos XV, XVI, y XVII, hombres como Leonardo de Vinci, Nicolás Copérnico, Galileo, Neper, William Harvey, Sir Francis Bacon y René Descartes, hicieron grandes operaciones al método científico, de tal forma que cuando se crearon los Estados Nacionales y surgió como fuerza el comercio internacional existía ya un método capaz de aplicarse a los datos económicos. Para el año 1532 empezaron a registrarse en Inglaterra las defunciones debido al temor que Enrique VII tenía por la peste. Más o menos por la misma época, en Francia la ley exigió a los clérigos registrar los bautismos, fallecimientos y matrimonios. Durante un brote de peste que apareció a fines de la década de 1500, el gobierno inglés comenzó a publicar estadística semanales de los decesos. Esa costumbre continuó muchos años, y en 1632 estos Bills of Mortality (Cuentas de Mortalidad) contenían los nacimientos y fallecimientos por sexo. En 1662, el capitán John Graunt usó documentos que abarcaban treinta años y efectuó predicciones sobre el número de personas que morirían de varias enfermedades y sobre las proporciones de nacimientos de varones y mujeres que cabría esperar. El trabajo de Graunt, condensado en su obra *Natural and Political Observations...Made upon the Bills of Mortality*, fue un esfuerzo innovador en el análisis estadístico. Por el año 1540 el alemán Sebastián Muster realizó una compilación estadística de los recursos nacionales, comprensiva de datos sobre organización política, instrucciones sociales, comercio y poderío militar.

Los eruditos del siglo XVII demostraron especial interés por la Estadística Demográfica como resultado de la especulación sobre si la población aumentaba, decrecía o permanecía estática. En los tiempos modernos tales métodos fueron resucitados por algunos reyes que necesitaban conocer las riquezas monetarias y el potencial humano de sus respectivos países. El primer empleo de los datos estadísticos para fines ajenos a la política tuvo lugar en 1691 y estuvo a cargo de Gaspar Neumann, un profesor alemán que vivía en Breslau. Este investigador se propuso destruir la antigua creencia popular de que en los años terminados en siete moría más gente que en los restantes, y para lograrlo hurgó pacientemente en los archivos parroquiales de la ciudad. Después de revisar miles de partidas de defunción pudo demostrar que en tales años no fallecían más personas que en los demás. Los procedimientos de Neumann fueron conocidos por el astrónomo inglés Halley, descubridor del cometa que lleva su nombre, quien los aplicó al estudio de la vida humana.

Durante el siglo XVII y principios del XVIII, matemáticos como Bernoulli, Francis Maseres, Lagrange y Laplace desarrollaron la teoría de probabilidades. No obstante durante cierto tiempo, la teoría de las

probabilidades limitó su aplicación a los juegos de azar y hasta el siglo XVIII no comenzó a aplicarse a los grandes problemas científicos. Godofredo Achenwall, profesor de la Universidad de Gotinga, acuñó en 1760 la palabra estadística, que extrajo del término italiano *statista* (estadista). Creía, y con sobrada razón, que los datos de la nueva ciencia serían el aliado más eficaz del gobernante consciente. La raíz remota de la palabra se halla, por otra parte, en el término latino *status*, que significa estado o situación; Esta etimología aumenta el valor intrínseco de la palabra, por cuanto la estadística revela el sentido cuantitativo de las más variadas situaciones. Jacques Quételet es quien aplica las Estadísticas a las ciencias sociales. Este interpretó la teoría de la probabilidad para su uso en las ciencias sociales y resolver la aplicación del principio de promedios y de la variabilidad a los fenómenos sociales. Quételet fue el primero en realizar la aplicación práctica de todo el método Estadístico, entonces conocido, a las diversas ramas de la ciencia. Entretanto, en el período del 1800 al 1820 se desarrollaron dos conceptos matemáticos fundamentales para la teoría Estadística; la teoría de los errores de observación, aportada por Laplace y Gauss; y la teoría de los mínimos cuadrados desarrollada por Laplace, Gauss y Legendre. A finales del siglo XIX, Sir Francis Gaston ideó el método conocido por Correlación, que tenía por objeto medir la influencia relativa de los factores sobre las variables. De aquí partió el desarrollo del coeficiente de correlación creado por Karl Pearson y otros cultivadores de la ciencia biométrica como J. Pease Norton, R. H. Hooker y G. Udny Yule, que efectuaron amplios estudios sobre la medida de las relaciones.

La historia de la estadística está resumida en tres grandes etapas o fases.

- **Fase 1: Los Censos:** Desde el momento en que se constituye una autoridad política, la idea de inventariar de una forma más o menos regular la población y las riquezas existentes en el territorio está ligada a la conciencia de soberanía y a los primeros esfuerzos administrativos.
- **Fase 2: De la Descripción de los Conjuntos a la Aritmética Política:** Las ideas mercantilistas extrañan una intensificación de este tipo de investigación. Colbert multiplica las encuestas sobre artículos manufacturados, el comercio y la población: los intendentes del Reino envían a París sus memorias. Vauban, más conocido por sus fortificaciones o su *Dime Royale*, que es la primera propuesta de un impuesto sobre los ingresos, se señala como el verdadero precursor de los sondeos. Más tarde, Bufón se preocupa de esos problemas antes de dedicarse a la historia natural. La escuela inglesa proporciona un nuevo progreso al superar la fase puramente descriptiva.
Sus tres principales representantes son Graunt, Petty y Halley. El penúltimo es autor de la famosa *Aritmética Política*. Chaptal, ministro del interior francés, publica en 1801 el primer censo general de población, desarrolla los estudios industriales, de las producciones y los cambios, haciéndose sistemáticos durante las dos terceras partes del siglo XIX.
- **Fase 3: Estadística y Cálculo de Probabilidades:** El cálculo de probabilidades se incorpora rápidamente como un instrumento de análisis extremadamente poderoso para el estudio de los fenómenos económicos y sociales y en general para el estudio de fenómenos cuyas causas son demasiadas complejas para conocerlos totalmente y hacer posible su análisis.

La Estadística para su mejor estudio se ha dividido en dos grandes ramas: **la Estadística Descriptiva y la Estadística Inferencial.**

- **Descriptiva:** consiste sobre todo en la presentación de datos en forma de tablas y gráficas. Esta comprende cualquier actividad relacionada con los datos y está diseñada para resumir o describir los mismos sin factores pertinentes adicionales; esto es, sin intentar inferir nada que vaya más allá de los datos, como tales.
- **Inferencial:** se deriva de muestras, de observaciones hechas sólo acerca de una parte de un conjunto numeroso de elementos y esto implica que su análisis requiere de generalizaciones que van más allá de los datos. Como consecuencia, la característica más importante del reciente crecimiento de la estadística ha sido un cambio en el énfasis de los métodos que describen a métodos que sirven para

hacer generalizaciones. La Estadística Inferencial investiga o analiza una población partiendo de una muestra tomada.

Estadística Inferencial

Los métodos básicos de la estadística inferencial son la estimación y el contraste de hipótesis, que juegan un papel fundamental en la investigación. Por tanto, algunos de los objetivos que se persiguen son:

- Calcular los parámetros de la distribución de medias o proporciones muestrales de tamaño n , extraídas de una población de media y varianza conocidas.
- Estimar la media o la proporción de una población a partir de la media o proporción muestral.
- Utilizar distintos tamaños muestrales para controlar la confianza y el error admitido.
- Contrastar los resultados obtenidos a partir de muestras.
- Visualizar gráficamente, mediante las respectivas curvas normales, las estimaciones realizadas.

En definitiva, la idea es, a partir de una población se extrae una muestra por algunos de los métodos existentes, con la que se generan datos numéricos que se van a utilizar para generar estadísticos con los que realizar estimaciones o contrastes poblacionales. Existen dos formas de estimar parámetros: la *estimación puntual* y la *estimación por intervalo de confianza*. En la primera se busca, con base en los datos muestrales, un único valor estimado para el parámetro. Para la segunda, se determina un intervalo dentro del cual se encuentra el valor del parámetro, con una probabilidad determinada.

Si el objetivo del tratamiento estadístico inferencial, es efectuar generalizaciones acerca de la estructura, composición o comportamiento de las poblaciones no observadas, a partir de una parte de la población, será necesario que la proporción de población examinada sea representativa del total. Por ello, la selección de la muestra requiere unos requisitos que lo garanticen, debe ser representativa y aleatoria.

Además, la cantidad de elementos que integran la muestra (el tamaño de la muestra) depende de múltiples factores, como el dinero y el tiempo disponibles para el estudio, la importancia del tema analizado, la confiabilidad que se espera de los resultados, las características propias del fenómeno analizado, etcétera.

Así, a partir de la muestra seleccionada se realizan algunos cálculos y se estima el valor de los parámetros de la población tales como la media, la varianza, la desviación estándar, o la forma de la distribución, etc.

El conjunto de los métodos que se utilizan para medir las características de la información, para resumir los valores individuales, y para analizar los datos a fin de extraerles el máximo de información, es lo que se llama *métodos estadísticos*. Los métodos de análisis para la información cuantitativa se pueden dividir en los siguientes seis pasos:

- Definición del problema.
- Recopilación de la información existente.
- Obtención de información original.
- Clasificación.
- Presentación.
- Análisis.

El centro de gravedad de la metodología estadística se empieza a desplazar técnicas de computación intensiva aplicadas a grandes masas de datos, y se empieza a considerar el método estadístico como un proceso iterativo de búsqueda del modelo ideal. Las aplicaciones en este periodo de la Estadística a la Economía

conducen a una disciplina con contenido propio: la Econometría. La investigación estadística en problemas militares durante la segunda guerra mundial y los nuevos métodos de programación matemática, dan lugar a la Investigación Operativa. El tratamiento de los datos de la investigación científica tiene varias etapas:

- En la etapa de recolección de datos del método científico, se define a la población de interés y se selecciona una muestra o conjunto de personas representativas de la misma, se realizan experimentos o se emplean instrumentos ya existentes o de nueva creación, para medir los atributos de interés necesarios para responder a las preguntas de investigación. Durante lo que es llamado trabajo de campo se obtienen los datos en crudo, es decir las respuestas directas de los sujetos uno por uno, se codifican (se les asignan valores a las respuestas), se capturan y se verifican para ser utilizados en las siguientes etapas.
- En la etapa de recuento, se organizan y ordenan los datos obtenidos de la muestra. Esta será descrita en la siguiente etapa utilizando la estadística descriptiva, todas las investigaciones utilizan estadística descriptiva, para conocer de manera organizada y resumida las características de la muestra.
- En la etapa de análisis se utilizan las pruebas estadísticas (estadística inferencial) y en la interpretación se acepta o rechaza la hipótesis nula.

Niveles de medición y tipos de variables

Para poder emplear el método estadístico en un estudio es necesario medir las variables.

- Medir: es asignar valores a las propiedades de los objetos bajo ciertas reglas, esas reglas son los niveles de medición.
- Cuantificar: es asignar valores a algo tomando un patrón de referencia. Por ejemplo, cuantificar es ver cuántos hombres y cuántas mujeres hay.

Variable: es una característica o propiedad que asume diferentes valores dentro de una población de interés y cuya variación es susceptible de medirse.

Las variables pueden clasificarse de acuerdo al tipo de valores que puede tomar como:

- **Discretas o categóricas** en las que los valores se relacionan a nombres, etiquetas o categorías, no existe un significado numérico directo.
- **Continuas** los valores tienen un correlato numérico directo, son continuos y susceptibles de fraccionarse y de poder utilizarse en operaciones aritméticas.
- **Dicotómica** sólo tienen dos valores posibles, la característica está ausente o presente.

En cuanto a una clasificación estadística, las variables pueden ser:

- **Aleatoria** Aquella en la cual desconocemos el valor porque fluctúa de acuerdo a un evento debido al azar.
- **Determinística** Aquella variable de la que se conoce el valor.
- **Independiente** aquellas variables que son manipuladas por el investigador. Define los grupos.
- **Dependiente** son mediciones que ocurren durante el experimento o tratamiento (resultado de la independiente), es la que se mide y compara entre los grupos.

En lo que tiene que ver con los **Niveles de Medición** tenemos distintos tipos de variable

- **Nominal:** Las propiedades de la medición nominal son:
 - * Exhaustiva: implica a todas las opciones.
 - * A los sujetos se les asignan categorías, por lo que son mutuamente excluyentes. Es decir, la variable está presente o no; tiene o no una característica.
- **Ordinal:** Las propiedades de la medición ordinal son:
 - * El nivel ordinal posee transitividad, por lo que se tiene la capacidad de identificar que es mejor o mayor que otra, en ese sentido se pueden establecer jerarquías.
 - * Las distancias entre un valor y otro no son iguales.
- **Intervalo:**
 - * El nivel de medición intervalar requiere distancias iguales entre cada valor. Por lo general utiliza datos cuantitativos. Por ejemplo: temperatura, atributos psicológicos (CI, nivel de autoestima, pruebas de conocimientos, etc.)
 - * Las unidades de calificación son equivalentes en todos los puntos de la escala. Una escala de intervalos implica: clasificación, magnitud y unidades de tamaños iguales (Brown, 2000).
 - * Se pueden hacer operaciones aritméticas.
 - * Cuando se le pide al sujeto que califique una situación del 0 al 10 puede tomarse como un nivel de medición de intervalo, siempre y cuando se incluya el 0.
- **Razón:**
 - * La escala empieza a partir del 0 absoluto, por lo tanto incluye sólo los números por su valor en sí, por lo que no pueden existir los números con signo negativo. Por ejemplo: Peso corporal en kg., edad en años, estatura en cm.

Definiciones adicionales

- **Variable:** Consideraciones que una variable son una característica o fenómeno que puede tomar distintos valores.
- **Dato:** Mediciones o cualidades que han sido recopiladas como resultado de observaciones.
- **Población:** Se considera el área de la cual son extraídos los datos. Es decir, es el conjunto de elementos o individuos que poseen una característica común y medible acerca de lo cual se desea información. Es también llamado Universo.
- **Muestra:** Es un subconjunto de la población, seleccionado de acuerdo a una regla o algún plan de muestreo.
- **Censo:** Recopilación de todos los datos (de interés para la investigación) de la población.
- **Estadística:** Es una función o fórmula que depende de los datos de la muestra (es variable).
- **Parámetro:** Característica medible de la población. Es un resumen numérico de alguna variable observada de la población. Los parámetros normales que se estudian son: *La media poblacional, Proporción.*
- **Estimador:** Un estimador de un parámetro es un estadístico que se emplea para conocer el parámetro desconocido.
- **Estadístico:** Es una función de los valores de la muestra. Es una variable aleatoria, cuyos valores dependen de la muestra seleccionada. Su distribución de probabilidad, se conoce como *Distribución muestral del estadístico.*

- **Estimación:** Este término indica que a partir de lo observado en una muestra (un resumen estadístico con las medidas que conocemos de Descriptiva) se extrapola o generaliza dicho resultado muestral a la población total, de modo que lo estimado es el valor generalizado a la población. Consiste en la búsqueda del valor de los parámetros poblacionales objeto de estudio. Puede ser puntual o por intervalo de confianza:

- * **Puntual:** cuando buscamos un valor concreto. Un estimador de un parámetro poblacional es una función de los datos muestrales. En pocas palabras, es una fórmula que depende de los valores obtenidos de una muestra, para realizar estimaciones. Lo que se pretende obtener es el valor exacto de un parámetro.

- * **Intervalo de confianza:** cuando determinamos un intervalo, dentro del cual se supone que va a estar el valor del parámetro que se busca con una cierta probabilidad. El intervalo de confianza está determinado por dos valores dentro de los cuales afirmamos que está el verdadero parámetro con cierta probabilidad. Son unos límites o margen de variabilidad que damos al valor estimado, para poder afirmar, bajo un criterio de probabilidad, que el verdadero valor no los rebasará.

Este intervalo contiene al parámetro estimado con una determinada certeza o nivel de confianza.

En la estimación por intervalos se usan los siguientes conceptos:

- **Variabilidad del parámetro:** Si no se conoce, puede obtenerse una aproximación en los datos o en un estudio piloto. También hay métodos para calcular el tamaño de la muestra que precinden de este aspecto. Habitualmente se usa como medida de esta variabilidad la desviación típica poblacional.
- **Error de la estimación:** Es una medida de su precisión que se corresponde con la amplitud del intervalo de confianza. Cuanta más precisión se desee en la estimación de un parámetro, más estrecho deberá ser el intervalo de confianza y, por tanto, menor el error, y más sujetos deberán incluirse en la muestra estudiada.
- **Nivel de confianza:** Es la probabilidad de que el verdadero valor del parámetro estimado en la población se sitúe en el intervalo de confianza obtenido. El nivel de confianza se denota por $1 - \alpha$
- **p-value:** También llamado nivel de significación. Es la probabilidad (en tanto por uno) de fallar en nuestra estimación, esto es, la diferencia entre la certeza (1) y el nivel de confianza $1 - \alpha$.
- **Valor crítico:** Se representa por $Z_{\alpha/2}$. Es el valor de la abscisa en una determinada distribución que deja a su derecha un área igual a $1/2$, siendo $1 - \alpha$ el nivel de confianza. Normalmente los valores críticos están tabulados o pueden calcularse en función de la distribución de la población.

Para un tamaño fijo de la muestra, los conceptos de error y nivel de confianza van relacionados. Si admitimos un error mayor, esto es, aumentamos el tamaño del intervalo de confianza, tenemos también una mayor probabilidad de éxito en nuestra estimación, es decir, un mayor nivel de confianza. Por tanto, un aspecto que debe de tenerse en cuenta es el tamaño muestral, ya que para disminuir el error que se comente habrá que aumentar el tamaño muestral. Esto se resolverá, para un intervalo de confianza cualquiera, despejando el tamaño de la muestra en cualquiera de las formulas de los intervalos de confianza que veremos a continuación, a partir del error máximo permitido. Los intervalos de confianza pueden ser unilaterales o bilaterales:

- **Contraste de Hipótesis:** Consiste en determinar si es aceptable, partiendo de datos muestrales, que la característica o el parámetro poblacional estudiado tome un determinado valor o esté dentro de unos determinados valores.
- **Nivel de Confianza:** Indica la proporción de veces que acertaríamos al afirmar que el parámetro está dentro del intervalo al seleccionar muchas muestras.

26.1.2 Muestreo:

Muestreo: Una muestra es representativa en la medida que es imagen de la población. En general, podemos decir que el tamaño de una muestra dependerá principalmente de: *Nivel de precisión deseado*, *Recursos disponibles*, *Tiempo involucrado en la investigación*. Además el plan de muestreo debe considerar *La población*, *Parámetros a medir*. Existe una gran cantidad de tipos de muestreo, en la práctica los más utilizados son los siguientes:

- **MUESTREO ALEATORIO SIMPLE:** Es un método de selección de n unidades extraídas de N , de tal manera que cada una de las posibles muestras tiene la misma probabilidad de ser escogida. (En la práctica, se enumeran las unidades de 1 a N , y a continuación se seleccionan n números aleatorios entre 1 y N , ya sea de tablas o de alguna urna con fichas numeradas).
- **MUESTREO ESTRATIFICADO ALEATORIO:** Se usa cuando la población está agrupada en pocos estratos, cada uno de ellos son muchas entidades. Este muestreo consiste en sacar una muestra aleatoria simple de cada uno de los estratos. (Generalmente, de tamaño proporcional al estrato).
- **MUESTREO SISTEMÁTICO:** Se utiliza cuando las unidades de la población están de alguna manera totalmente ordenadas. Para seleccionar una muestra de n unidades, se divide la población en n subpoblaciones de tamaño $K = N/n$ y se toma al azar una unidad de la K primeras y de ahí en adelante cada K -ésima unidad.
- **MUESTREO POR CONGLOMERADO:** Se emplea cuando la población está dividida en grupos o conglomerados pequeños. Consiste en obtener una muestra aleatoria simple de conglomerados y luego CENSAR cada uno de éstos.
- **MUESTREO EN DOS ETAPAS (Bietápico):** En este caso la muestra se toma en dos pasos:
 - * Seleccionar una muestra de unidades primarias, y
 - * Seleccionar una muestra de elementos a partir de cada unidad primaria escogida.
 - * *Observación:* En la realidad es posible encontrarse con situaciones en las cuales no es posible aplicar libremente un tipo de muestreo, incluso estaremos obligados a mezclarlas en ocasiones.

26.1.3 Errores Estadísticos Comunes

El propósito de esta sección es solamente indicar los malos usos comunes de datos estadísticos, sin incluir el uso de métodos estadísticos complicados. Un estudiante debería estar alerta en relación con estos malos usos y debería hacer un gran esfuerzo para evitarlos a fin de ser un verdadero estadístico.

Datos estadísticos inadecuados: Los datos estadísticos son usados como la materia prima para un estudio estadístico. Cuando los datos son inadecuados, la conclusión extraída del estudio de los datos se vuelve obviamente inválida. Por ejemplo, supongamos que deseamos encontrar el ingreso familiar típico del año pasado en la ciudad Y de 50,000 familias y tenemos una muestra consistente del ingreso de solamente tres familias: 1 millón, 2 millones y no ingreso. Si sumamos el ingreso de las tres familias y dividimos el total por 3, obtenemos un promedio de 1 millón. Entonces, extraemos una conclusión basada en la muestra de que el ingreso familiar promedio durante el año pasado en la ciudad fue de 1 millón. Es obvio que la conclusión es falsa, puesto que las cifras son extremas y el tamaño de la muestra es demasiado pequeño; por lo tanto la muestra no es representativa.

Hay muchas otras clases de datos inadecuados. Por ejemplo, algunos datos son respuestas inexactas de una encuesta, porque las preguntas usadas en la misma son vagas o engañosas, algunos datos son toscas estimaciones porque no hay disponibles datos exactos o es demasiado costosa su obtención, y algunos datos son irrelevantes en un problema dado, porque el estudio estadístico no está bien planeado. Al momento de recopilar los datos que serán procesados se es susceptible de cometer errores así como

durante los cálculos de los mismos. No obstante, hay otros errores que no tienen nada que ver con la digitación y que no son tan fácilmente identificables. Algunos de éstos errores son:

- **Sesgo:** Es imposible ser completamente objetivo o no tener ideas preconcebidas antes de comenzar a estudiar un problema, y existen muchas maneras en que una perspectiva o estado mental pueda influir en la recopilación y en el análisis de la información. En estos casos se dice que hay un sesgo cuando el individuo da mayor peso a los datos que apoyan su opinión que a aquellos que la contradicen. Un caso extremo de sesgo sería la situación donde primero se toma una decisión y después se utiliza el análisis estadístico para justificar la decisión ya tomada.
- **Datos No Comparables:** el establecer comparaciones es una de las partes más importantes del análisis estadístico, pero es extremadamente importante que tales comparaciones se hagan entre datos que sean comparables.
- **Proyección descuidada de tendencias:** la proyección simplista de tendencias pasadas hacia el futuro es uno de los errores que más ha desacreditado el uso del análisis estadístico.
- **Muestreo Incorrecto:** en la mayoría de los estudios sucede que el volumen de información disponible es tan inmenso que se hace necesario estudiar muestras, para derivar conclusiones acerca de la población a que pertenece la muestra. Si la muestra se selecciona correctamente, tendrá básicamente las mismas propiedades que la población de la cual fue extraída; pero si el muestreo se realiza incorrectamente, entonces puede suceder que los resultados no signifiquen nada.

Sesgo significa que un usuario dé los datos perjudicialmente de más énfasis a los hechos, los cuales son empleados para mantener su predeterminada posición u opinión. Los estadísticos son frecuentemente degradados por lemas tales como: *Hay tres clases de mentiras: mentiras, mentiras reprobables y estadística, y Las cifras no mienten, pero los mentirosos piensan.* Hay dos clases de sesgos: conscientes e inconscientes. Ambos son comunes en el análisis estadístico. Hay numerosos ejemplos de sesgos conscientes. Un anunciante frecuentemente usa la estadística para probar que su producto es muy superior al producto de su competidor. Un político prefiere usar la estadística para sostener su punto de vista. Gerentes y líderes de trabajadores pueden simultáneamente situar sus respectivas cifras estadísticas sobre la misma tabla de trato para mostrar que sus rechazos o peticiones son justificadas. Es casi imposible que un sesgo inconsciente esté completamente ausente en un trabajo estadístico. En lo que respecta al ser humano, es difícil obtener una actitud completamente objetiva al abordar un problema, aun cuando un científico debería tener una mente abierta. Un estadístico debería estar enterado del hecho de que su interpretación de los resultados del análisis estadístico está influenciado por su propia experiencia, conocimiento y antecedentes con relación al problema dado.

26.2 Fundamentos

26.2.1 Pruebas de Hipótesis

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado.
- El interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.

- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

Una prueba de hipótesis está formada por cinco partes:

- La hipótesis nula, denotada por H_0 .
- La hipótesis alternativa, denotada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 24. Las dos hipótesis en competencia son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones:

- Rechazar H_0 y concluir que H_1 es verdadera.
- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 11. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19 usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 12. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas:

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 25. El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadístico de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**. Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**. Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente. Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 26. Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

26.2.2 Muestras grandes: una media poblacional

Definición 27. El **valor de p** (*p-value*) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.

Nota 24. Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.

Nota 25. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 28. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados:

p	H_0	Significativa
$p \leq 0.01$	rechazada	Result. altamente significativos y en contra de H_0
$p \leq 0.05$	rechazada	Result. Estadísticamente significativos y en contra de H_0
$p \leq 0.10$	rechazada	Result. posiblemente significativos con Tendencia estadística y en contra de H_0
$p > 0.10$	no rechazada	Result. estadísticamente no significativos y no rechazar H_0

Nota 26. Para determinar el valor de p , encontrar el área en la cola después del estadístico de prueba. Si la prueba es de una cola, este es el valor de p . Si es de dos colas, éste valor encontrado es la mitad del valor de p . Rechazar H_0 cuando el valor de $p < \alpha$.

Hay dos tipos de errores al realizar una prueba de hipótesis

	H_0 es Verdadera	H_0 es Falsa
Rechazar H_0	Error tipo I	✓
Aceptar H_0	✓	Error tipo II

Definición 29. La probabilidad de cometer el error tipo II se define por β donde

$$\begin{aligned}\beta &= P\{\text{error tipo II}\} = P\{\text{Aceptar equivocadamente } H_0\} \\ &= P\{\text{Aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\}\end{aligned}$$

Nota 27. Cuando H_0 es falsa y H_1 es verdadera, no siempre es posible especificar un valor exacto de μ , sino más bien un rango de posibles valores. En lugar de arriesgarse a tomar una decisión incorrecta, es mejor concluir que no hay evidencia suficiente para rechazar H_0 , es decir en lugar de aceptar H_0 , no rechazar H_0 .

La bondad de una prueba estadística se mide por el tamaño de α y β , ambas deben de ser pequeñas. Una manera muy efectiva de medir la potencia de la prueba es calculando el complemento del error tipo II:

$$\begin{aligned} 1 - \beta &= P\{\text{Rechazar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{\text{Rechazar } H_0 \text{ cuando } H_1 \text{ es verdadera}\} \end{aligned}$$

Definición 30. La **potencia de la prueba**, $1 - \beta$, mide la capacidad de que la prueba funciona como se necesita.

Ejemplo 13. La producción diaria de una planta química local ha promediado 880 toneladas en los últimos años. A la gerente de control de calidad le gustaría saber si este promedio ha cambiado en meses recientes. Ella selecciona al azar 50 días de la base de datos computarizada y calcula el promedio y la desviación estándar de las $n = 50$ producciones como $\bar{x} = 871$ toneladas y $s = 21$ toneladas, respectivamente. Pruebe la hipótesis apropiada usando $\alpha = 0.05$.

La hipótesis nula apropiada es:

$$\begin{aligned} H_0 &: \mu = 880 \\ &\text{y la hipótesis alternativa } H_1 \text{ es} \\ H_1 &: \mu \neq 880 \end{aligned}$$

el estimador puntual para μ es \bar{x} , entonces el estadístico de prueba es

$$\begin{aligned} z &= \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \\ &= \frac{871 - 880}{21/\sqrt{50}} = -3.03 \end{aligned}$$

Para esta prueba de dos colas, hay que determinar los dos valores de $z_{\alpha/2}$, es decir, $z_{\alpha/2} = \pm 1.96$, como $z > z_{\alpha/2}$, z cae en la zona de rechazo, por lo tanto la gerente puede rechazar la hipótesis nula y concluir que el promedio efectivamente ha cambiado. La probabilidad de rechazar H_0 cuando esta es verdadera es de 0.05. Recordemos que el valor observado del estadístico de prueba es $z = -3.03$, la región de rechazo más pequeña que puede usarse y todavía seguir rechazando H_0 es $|z| > 3.03$, entonces $p = 2(0.012) = 0.0024$, que a su vez es menor que el nivel de significancia α asignado inicialmente, y además los resultados son **altamente significativos**. Finalmente determinemos la potencia de la prueba cuando μ en realidad es igual a 870 toneladas.

Recordar que la región de aceptación está entre -1.96 y 1.96 , para $\mu = 880$, equivalentemente

$$874.18 < \bar{x} < 885.82$$

β es la probabilidad de aceptar H_0 cuando $\mu = 870$, calculemos los valores de z correspondientes a 874.18 y 885.82. Entonces

$$\begin{aligned} z_1 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{874.18 - 870}{21/\sqrt{50}} = 1.41 \\ z_2 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{885.82 - 870}{21/\sqrt{50}} = 5.33 \end{aligned}$$

por lo tanto

$$\begin{aligned}\beta &= P\{\text{aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\} = P\{874.18 < \mu < 885.82 \text{ cuando } \mu = 870\} \\ &= P\{1.41 < z < 5.33\} = P\{1.41 < z\} = 1 - 0.9207 = 0.0793\end{aligned}$$

entonces, la potencia de la prueba es

$$1 - \beta = 1 - 0.0793 = 0.9207$$

que es la probabilidad de rechazar correctamente H_0 cuando H_0 es falsa.

Prueba de hipótesis para la diferencia entre dos medias poblacionales

El estadístico que resume la información muestral respecto a la diferencia en medias poblacionales $(\mu_1 - \mu_2)$ es la diferencia de las medias muestrales $(\bar{x}_1 - \bar{x}_2)$, por tanto al probar la diferencia entre las medias muestrales se verifica que la diferencia real entre las medias poblacionales difiere de un valor especificado, $(\mu_1 - \mu_2) = D_0$, se puede usar el error estándar de $(\bar{x}_1 - \bar{x}_2)$, es decir

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

cuyo estimador está dado por

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

El procedimiento para muestras grandes es:

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$, donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.
- 2) **Hipótesis Alternativa**

Prueba de una Cola	Prueba de dos colas
$H_1 : (\mu_1 - \mu_2) > D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$
$H_1 : (\mu_1 - \mu_2) < D_0$	

- 3) Estadístico de prueba:

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_0$	
$z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cuando $p < \alpha$	

Ejemplo 14. Para determinar si ser propietario de un automóvil afecta el rendimiento académico de un estudiante, se tomaron dos muestras aleatorias de 100 estudiantes varones. El promedio de calificaciones para los $n_1 = 100$ no propietarios de un auto tuvieron un promedio y varianza de $\bar{x}_1 = 2.7$ y $s_1^2 = 0.36$, respectivamente, mientras que para la segunda muestra con $n_2 = 100$ propietarios de un auto, se tiene $\bar{x}_2 = 2.54$ y $s_2^2 = 0.4$. Los datos presentan suficiente evidencia para indicar una diferencia en la media en el rendimiento académico entre propietarios y no propietarios de un automóvil? Hacer pruebas para $\alpha = 0.01, 0.05$ y $\alpha = 0.1$.

- Solución utilizando la técnica de regiones de rechazo: realizando las operaciones $z = 1.84$, determinar si excede los valores de $z_{\alpha/2}$.
- Solución utilizando el p-value: Calcular el valor de p , la probabilidad de que z sea mayor que $z = 1.84$ o menor que $z = -1.84$, se tiene que $p = 0.0658$.
- Si el intervalo de confianza que se construye contiene el valor del parámetro especificado por H_0 , entonces ese valor es uno de los posibles valores del parámetro y H_0 no debe ser rechazada.
- Si el valor hipotético se encuentra fuera de los límites de confianza, la hipótesis nula es rechazada al nivel de significancia α .

Prueba de Hipótesis para una Proporción Binomial

Para una muestra aleatoria de n intentos idénticos, de una población binomial, la proporción muestral \hat{p} tiene una distribución aproximadamente normal cuando n es grande, con media p y error estándar

$$SE = \sqrt{\frac{pq}{n}}.$$

La prueba de hipótesis de la forma

$$\begin{aligned} H_0 &: p = p_0 \\ H_1 &: p > p_0, \text{ o } p < p_0 \text{ o } p \neq p_0 \end{aligned}$$

El estadístico de prueba se construye con el mejor estimador de la proporción verdadera, \hat{p} , con el estadístico de prueba z , que se distribuye normal estándar.

El procedimiento es

- 1) Hipótesis nula: $H_0 : p = p_0$
- 2) Hipótesis alternativa

Prueba de una Cola	Prueba de dos colas
$H_1 : p > p_0$	$p \neq p_0$
$H_1 : p < p_0$	

- 3) Estadístico de prueba:

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{pq}{n}}}, \hat{p} = \frac{x}{n}$$

donde x es el número de éxitos en n intentos binomiales.

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_0$	
$z < -z_\alpha$ cuando $H_1 : p < p_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cundo $p < \alpha$	

Prueba de Hipótesis diferencia entre dos Proporciones Binomiales

Cuando se tienen dos muestras aleatorias independientes de dos poblaciones binomiales, el objetivo del experimento puede ser la diferencia $(p_1 - p_2)$ en las proporciones de individuos u objetos que poseen una característica específica en las dos poblaciones. En este caso se pueden utilizar los estimadores de las dos proporciones $(\hat{p}_1 - \hat{p}_2)$ con error estándar dado por

$$SE = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}},$$

considerando el estadístico z con un nivel de significancia $(1 - \alpha) 100\%$

La hipótesis nula a probarse es de la forma

$H_0: p_1 = p_2$ o equivalentemente $(p_1 - p_2) = 0$, contra una hipótesis alternativa H_1 de una o dos colas.

Para estimar el error estándar del estadístico z , se debe de utilizar el hecho de que suponiendo que H_0 es verdadera, las dos proporciones son iguales a algún valor común, p . Para obtener el mejor estimador de p es

$$p = \frac{\text{número total de éxitos}}{\text{Número total de pruebas}} = \frac{x_1 + x_2}{n_1 + n_2}.$$

1) **Hipótesis Nula:** $H_0 : (p_1 - p_2) = 0$

2) **Hipótesis Alternativa:** $H_1 :$

Prueba de una Cola	Prueba de dos colas
$H_1 : (p_1 - p_2) > 0$	$H_1 : (p_1 - p_2) \neq 0$
$H_1 : (p_1 - p_2) < 0$	

3) Estadístico de prueba:

$$z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{pq}{n_1} + \frac{pq}{n_2}}},$$

donde $\hat{p}_1 = x_1/n_1$ y $\hat{p}_2 = x_2/n_2$, dado que el valor común para p_1 y p_2 es p , entonces $\hat{p} = \frac{x_1+x_2}{n_1+n_2}$ y por tanto el estadístico de prueba es

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}.$$

4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_\alpha$	
$z < -z_\alpha$ cuando $H_1 : p < p_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cundo $p < \alpha$	

26.2.3 Muestras Pequeñas

1) **Hipótesis Nula:** $H_0 : \mu = \mu_0$,

2) **Hipótesis Alternativa:** $H_1 :$

Prueba de una Cola	Prueba de dos colas
$H_1 : \mu > \mu_0$	$H_1 : \mu \neq \mu_0$
$H_1 : \mu < \mu_0$	

3) Estadístico de prueba:

$$t = \frac{\bar{x} - \mu_0}{\sqrt{\frac{s^2}{n}}},$$

4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$t > t_\alpha$	
$t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Diferencia entre dos medias poblacionales: MAI

Cuando los tamaños de muestra son pequeños, no se puede asegurar que las medias muestrales sean normales, pero si las poblaciones originales son normales, entonces la distribución muestral de la diferencia de las medias muestrales, $(\bar{x}_1 - \bar{x}_2)$, será normal con media $(\mu_1 - \mu_2)$ y error estándar

$$ES = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}.$$

1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

2) **Hipótesis Alternativa**

Prueba de una Cola	Prueba de dos colas
$H_1 : (\mu_1 - \mu_2) > D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$
$H_1 : (\mu_1 - \mu_2) < D_0$	

3) Estadístico de prueba:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

donde

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

4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_\alpha$	
$z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

Diferencia entre dos medias poblacionales: Diferencias Pareadas

1) **Hipótesis Nula:** $H_0 : \mu_d = 0$

2) **Hipótesis Alternativa:** $H_1 : \mu_d$

Prueba de una Cola	Prueba de dos colas
$H_1 : \mu_d > 0$	$H_1 : \mu_d \neq 0$
$H_1 : \mu_d < 0$	

3) Estadístico de prueba:

$$t = \frac{\bar{d}}{\sqrt{\frac{s_d^2}{n}}}$$

donde n es el número de diferencias pareadas, \bar{d} es la media de las diferencias muestrales, y s_d es la desviación estándar de las diferencias muestrales.

4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$t > t_\alpha$	
$t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$
cuando $p < \alpha$	

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

Inferencias con respecto a la Varianza Poblacional

1) **Hipótesis Nula:** $H_0 : \sigma^2 = \sigma_0^2$

2) **Hipótesis Alternativa:** H_1

Prueba de una Cola	Prueba de dos colas
$H_1 : \sigma^2 > \sigma_0^2$	$H_1 : \sigma^2 \neq \sigma_0^2$
$H_1 : \sigma^2 < \sigma_0^2$	

3) Estadístico de prueba:

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2},$$

4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$\chi^2 > \chi_\alpha^2$	
$\chi^2 < \chi_{(1-\alpha)}^2$ cuando $H_1 : \chi^2 < \chi_0^2$	$\chi^2 > \chi_{\alpha/2}^2$ o $\chi^2 < \chi_{(1-\alpha/2)}^2$
cuando $p < \alpha$	

Los valores críticos de χ^2 , están basados en $(n_1 +)$ grados de libertad.

Comparación de dos varianzas poblacionales

1) **Hipótesis Nula** $H_0 : (\sigma_1^2 - \sigma_2^2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

2) **Hipótesis Alternativa**

Prueba de una Cola	Prueba de dos colas
$H_1 : (\sigma_1^2 - \sigma_2^2) > D_0$	$H_1 : (\sigma_1^2 - \sigma_2^2) \neq D_0$
$H_1 : (\sigma_1^2 - \sigma_2^2) < D_0$	

3) Estadístico de prueba:

$$F = \frac{s_1^2}{s_2^2}$$

donde s_1^2 es la varianza muestral más grande.

4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$F > F_\alpha$	$F > F_{\alpha/2}$
cuando $p < \alpha$	

26.2.4 Estimación por intervalos

Recordemos que S^2 es un estimador insesgado de σ^2 , entonces se tiene la siguiente definición

Definición 31. Sean $\hat{\theta}_1$ y $\hat{\theta}_2$ dos estimadores insesgados de θ , parámetro poblacional. Si $\sigma_{\hat{\theta}_1}^2 < \sigma_{\hat{\theta}_2}^2$, decimos que $\hat{\theta}_1$ es un estimador más eficaz de θ que $\hat{\theta}_2$.

Algunas observaciones que es preciso realizar

Nota 28. a) Para poblaciones normales, \bar{X} y \tilde{X} son estimadores insesgados de μ , pero con $\sigma_{\bar{X}}^2 < \sigma_{\tilde{X}}^2$.

b) Para las estimaciones por intervalos de θ , un intervalo de la forma $\hat{\theta}_L < \theta < \hat{\theta}_U$, $\hat{\theta}_L$ y $\hat{\theta}_U$ dependen del valor de $\hat{\theta}$.

c) Para $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$, si $n \rightarrow \infty$, entonces $\hat{\theta} \rightarrow \mu$.

Nota 29. Para $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$, si $n \rightarrow \infty$,

d) Para $\hat{\theta}$ se determinan $\hat{\theta}_L$ y $\hat{\theta}_U$ de modo tal que

$$P\left\{\hat{\theta}_L < \hat{\theta} < \hat{\theta}_U\right\} = 1 - \alpha,$$

con $\alpha \in (0, 1)$. Es decir, $\theta \in (\hat{\theta}_L, \hat{\theta}_U)$ es un intervalo de confianza del $100(1 - \alpha)\%$.

e) De acuerdo con el TLC se espera que la distribución muestral de \bar{X} se distribuye aproximadamente normal con media $\mu_{\bar{X}} = \mu$ y desviación estándar $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$.

Para $Z_{\alpha/2}$ se tiene $P\{-Z_{\alpha/2} < Z < Z_{\alpha/2}\} = 1 - \alpha$, donde $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$. Entonces

$$P\left\{-Z_{\alpha/2} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < Z_{\alpha/2}\right\} = 1 - \alpha,$$

es equivalente a

$$P\left\{\bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right\} = 1 - \alpha.$$

f) Si \bar{X} es la media muestral de una muestra de tamaño n de una población con varianza conocida σ^2 , el intervalo de confianza de $100(1 - \alpha)\%$ para μ es

$$\mu \in \left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$$

- g) Para muestras pequeñas de poblaciones no normales, no se puede esperar que el grado de confianza sea preciso.
- h) Para $n \geq 30$, con distribución de forma no muy sesgada, se pueden tener buenos resultados.

Teorema 3. Si \bar{X} es un estimador de μ , podemos tener $100(1 - \alpha)\%$ de confianza en que el error no excederá a $z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$, error entre \bar{X} y μ .

Teorema 4. Si \bar{X} es un estimador de μ , podemos tener $100(1 - \alpha)\%$ de confianza en que el error no excederá una cantidad e cuando el tamaño de la muestra es

$$n = \left(\frac{z_{\alpha/2} \sigma}{e} \right)^2.$$

Nota 30. Para intervalos unilaterales

$$P \left\{ \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < Z_{\alpha} \right\} = 1 - \alpha$$

equivalentemente

$$P \left\{ \mu < \bar{X} + Z_{\alpha} \frac{\sigma}{\sqrt{n}} \right\} = 1 - \alpha.$$

Si \bar{X} es la media de una muestra aleatoria de tamaño n a partir de una población con varianza σ^2 , los límites de confianza unilaterales del $100(1 - \alpha)\%$ de confianza para μ están dados por

- a) Límite unilateral superior: $\bar{x} + z_{\alpha} \frac{\sigma}{\sqrt{n}}$
- b) Límite unilateral inferior: $\bar{x} - z_{\alpha} \frac{\sigma}{\sqrt{n}}$
- c) Para σ desconocida recordar que $T = \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_{n-1}$, donde s es la desviación estándar de la muestra. Entonces

$$P \left\{ -t_{\alpha/2} < T < t_{\alpha/2} \right\} = 1 - \alpha, \text{ equivalentemente } P \left\{ \bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}} \right\} = 1 - \alpha.$$

- d) Un intervalo de confianza del $100(1 - \alpha)\%$ de confianza para μ , σ^2 desconocida y población normal es $\mu \in \left(\bar{x} - t_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}} \right)$, donde $t_{\alpha/2}$ es una t -student con $\nu = n - 1$ grados de libertad.
- e) Los límites unilaterales para μ con σ desconocida son $\bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}}$ y $\bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}$.
- f) Cuando la población no es normal, σ desconocida y $n \geq 30$, σ se puede reemplazar por s para obtener el intervalo de confianza para muestras grandes:

$$\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}.$$

- g) El estimador de \bar{X} de μ , σ desconocida, la varianza de $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$, el error estándar de \bar{X} es σ/\sqrt{n} .
- h) Si σ es desconocida y la población es normal, $s \rightarrow \sigma$ y se incluye el error estándar s/\sqrt{n} , entonces

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}.$$

Intervalos de confianza sobre la varianza

Supongamos que X se distribuye normal (μ, σ^2) , desconocidas. Sea X_1, X_2, \dots, X_n muestra aleatoria de tamaño n , s^2 la varianza muestral.

Se sabe que $X^2 = \frac{(n-1)s^2}{\sigma^2}$ se distribuye χ_{n-1}^2 grados de libertad. Su intervalo de confianza es

$$\begin{aligned} P \left\{ \chi_{1-\frac{\alpha}{2}, n-1}^2 \leq \chi^2 \leq \chi_{\frac{\alpha}{2}, n-1}^2 \right\} &= 1 - \alpha \\ P \left\{ \chi_{1-\frac{\alpha}{2}, n-1}^2 \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi_{\frac{\alpha}{2}, n-1}^2 \right\} &= 1 - \alpha \\ P \left\{ \frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right\} &= 1 - \alpha, \end{aligned}$$

es decir

$$\sigma^2 \in \left[\frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2}, \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right],$$

los intervalos unilaterales son

$$\sigma^2 \in \left[\frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2}, \infty \right),$$

y

$$\sigma^2 \in \left(-\infty, \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right].$$

Intervalos de confianza para proporciones

Supongamos que se tienen una muestra de tamaño n de una población grande pero finita, y supongamos que X , $X \leq n$, pertenecen a la clase de interés, entonces

$$\hat{p} = \frac{\bar{X}}{n},$$

es el estimador puntual de la proporción de la población que pertenece a dicha clase. n y p son los parámetros de la distribución binomial, entonces $\hat{p} \sim N \left(p, \frac{p(1-p)}{n} \right)$ aproximadamente si p es distinto de 0 y 1; o si n es suficientemente grande. Entonces

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim N(0, 1), \text{ aproximadamente.}$$

Entonces

$$1 - \alpha = P \left\{ -z_{\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \leq z_{\alpha/2} \right\} = P \left\{ \hat{p} - z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \right\}$$

con $\sqrt{\frac{p(1-p)}{n}}$ error estándar del estimador puntual p . Una solución para determinar el intervalo de confianza del parámetro p (desconocido) es

$$1 - \alpha = P \left\{ \hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right\}$$

entonces los intervalos de confianza, tanto unilaterales como de dos colas son:

- a) $p \in \left(\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right),$
b) $p \in \left(-\infty, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right),$
c) $p \in \left(\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \infty \right);$

para minimizar el error estándar, se propone que el tamaño de la muestra sea

$$n = \left(\frac{z_{\alpha/2}}{E} \right)^2 p(1-p)$$

, donde

$$E = |p - \hat{p}|.$$

Intervalos de confianza para dos muestras: Varianzas conocidas

Sean X_1 y X_2 variables aleatorias independientes. X_1 con media desconocida μ_1 y varianza conocida σ_1^2 ; y X_2 con media desconocida μ_2 y varianza conocida σ_2^2 . Se busca encontrar un intervalo de confianza de $100(1-\alpha)\%$ de la diferencia entre medias μ_1 y μ_2 . Sean $X_{11}, X_{12}, \dots, X_{1n_1}$ muestra aleatoria de n_1 observaciones de X_1 , y sean $X_{21}, X_{22}, \dots, X_{2n_2}$ muestra aleatoria de n_2 observaciones de X_2 .

Sean \bar{X}_1 y \bar{X}_2 , medias muestrales, entonces el estadístico

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1),$$

si X_1 y X_2 son normales o aproximadamente normales si se aplican las condiciones del Teorema de Límite Central respectivamente. Entonces se tiene

$$\begin{aligned} 1 - \alpha &= P\{-Z_{\alpha/2} \leq Z \leq Z_{\alpha/2}\} = P\left\{-Z_{\alpha/2} \leq \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \leq Z_{\alpha/2}\right\} \\ &= P\left\{(\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \leq \mu_1 - \mu_2 \leq (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}\right\}. \end{aligned}$$

Entonces los intervalos de confianza unilaterales y de dos colas al $(1-\alpha)\%$ de confianza son

- a) $\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right],$
b) $\mu_1 - \mu_2 \in \left[-\infty, (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right],$
c) $\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \infty \right].$

Nota 31. Si σ_1 y σ_2 son conocidas, o por lo menos se conoce una aproximación, y los tamaños de las muestras n_1 y n_2 son iguales, $n_1 = n_2 = n$, se puede determinar el tamaño de la muestra para que el error al estimar $\mu_1 - \mu_2$ usando $\bar{X}_1 - \bar{X}_2$ sea menor que E (valor del error deseado) al $(1-\alpha)\%$ de confianza. El tamaño n de la muestra requerido para cada muestra es

$$n = \left(\frac{Z_{\alpha/2}}{E} \right)^2 (\sigma_1^2 + \sigma_2^2).$$

Intervalos de confianza para dos muestras: Varianzas desconocidas e iguales

- a) Si $n_1, n_2 \geq 30$ se pueden utilizar los intervalos de la distribución normal para varianza conocida
- b) Si n_1, n_2 son muestras pequeñas, supongase que las poblaciones para X_1 y X_2 son normales con varianzas desconocidas y con base en el intervalo de confianza para distribuciones t -student

Supongamos que X_1 es una variable aleatoria con media μ_1 y varianza σ_1^2 , X_2 es una variable aleatoria con media μ_2 y varianza σ_2^2 . Todos los parámetros son desconocidos. Sin embargo supóngase que es razonable considerar que $\sigma_1^2 = \sigma_2^2 = \sigma^2$.

Nuevamente sean X_1 y X_2 variables aleatorias independientes. X_1 con media desconocida μ_1 y varianza muestral S_1^2 ; y X_2 con media desconocida μ_2 y varianza muestral S_2^2 . Dado que S_1^2 y S_2^2 son estimadores de σ_1^2 , se propone el estimador S de σ^2 como

$$S_p^2 = \frac{(n_1 - 1) S_1^2 + (n_2 - 1) S_2^2}{n_1 + n_2 - 2},$$

entonces, el estadístico para $\mu_1 - \mu_2$ es

$$t_\nu = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$$

donde t_ν es una t de student con $\nu = n_1 + n_2 - 2$ grados de libertad.

Por lo tanto

$$\begin{aligned} 1 - \alpha &= P \{ -t_{\alpha/2, \nu} \leq t \leq t_{\alpha/2, \nu} \} \\ &= P \left\{ (\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \leq \right. \\ &\quad \left. t \leq (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right\}, \end{aligned}$$

luego, los intervalos de confianza del $(1 - \alpha) \%$ para $\mu_1 - \mu_2$ son

- a) $\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right].$
- b) $\mu_1 - \mu_2 \in \left[-\infty, (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right].$
- c) $\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \infty \right].$

Intervalos de confianza para dos muestras: Varianzas desconocidas diferentes

Si no se tiene certeza de que $\sigma_1^2 = \sigma_2^2$, se propone el estadístico

$$t^* = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}},$$

que se distribuye t -student con ν grados de libertad, donde

$$\nu = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\frac{S_1^2/n_1}{n_1+1} + \frac{S_2^2/n_2}{n_2+1}} - 2.$$

Entonces el intervalo de confianza de aproximadamente el $100(1 - \alpha)\%$ para $\mu_1 - \mu_2$ con $\sigma_1^2 \neq \sigma_2^2$ es

$$\mu_1 - \mu_2 \in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}, (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \right].$$

Intervalos de confianza para razón de Varianzas

Supongamos que se toman dos muestras aleatorias independientes de las dos poblaciones de interés. Sean X_1 y X_2 variables normales independientes con medias desconocidas μ_1 y μ_2 y varianzas desconocidas σ_1^2 y σ_2^2 respectivamente. Se busca un intervalo de confianza de $100(1 - \alpha)\%$ para σ_1^2/σ_2^2 . Supongamos n_1 y n_2 muestras aleatorias de X_1 y X_2 y sean S_1^2 y S_2^2 varianzas muestrales. Se sabe que

$$F = \frac{S_2^2/\sigma_2^2}{S_1^2/\sigma_1^2},$$

se distribuye F con $n_2 - 1$ y $n_1 - 1$ grados de libertad.

Por lo tanto

$$P\left\{F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq F \leq F_{\frac{\alpha}{2}, n_2-1, n_1-1}\right\} = 1 - \alpha,$$

$$P\left\{F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq \frac{S_2^2/\sigma_2^2}{S_1^2/\sigma_1^2} \leq F_{\frac{\alpha}{2}, n_2-1, n_1-1}\right\} = 1 - \alpha,$$

luego entonces

$$P\left\{\frac{S_1^2}{S_2^2} F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{S_1^2}{S_2^2} F_{\frac{\alpha}{2}, n_2-1, n_1-1}\right\} = 1 - \alpha.$$

en consecuencia

$$\frac{\sigma_1^2}{\sigma_2^2} \in \left[\frac{S_1^2}{S_2^2} F_{1-\frac{\alpha}{2}, n_2-1, n_1-1}, \frac{S_1^2}{S_2^2} F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right],$$

donde

$$F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} = \frac{1}{F_{\frac{\alpha}{2}, n_2-1, n_1-1}}.$$

Intervalos de confianza para diferencia de proporciones

Sean dos proporciones de interés p_1 y p_2 . Se busca un intervalo para $p_1 - p_2$ al $100(1 - \alpha)\%$. Sean dos muestras independientes de tamaño n_1 y n_2 de poblaciones infinitas de modo que X_1 y X_2 variables aleatorias binomiales independientes con parámetros (n_1, p_1) y (n_2, p_2) . X_1 y X_2 son el número de observaciones que pertenecen a la clase de interés correspondientes. Entonces $\hat{p}_1 = \frac{X_1}{n_1}$ y $\hat{p}_2 = \frac{X_2}{n_2}$ son estimadores de p_1 y p_2 respectivamente. Supongamos que se cumple la aproximación normal a la binomial, entonces

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n_1} - \frac{p_2(1-p_2)}{n_2}}} \sim N(0, 1) \text{ aproximadamente}$$

por tanto

$$(\hat{p}_1 - \hat{p}_2) - Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \leq p_1 - p_2 \leq (\hat{p}_1 - \hat{p}_2) + Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

26.3 Bases

26.3.1 Análisis de Regresión Lineal (RL)

Nota 32. — En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.
- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

26.3.2 Análisis de Regresión Lineal (RL)

Nota 33. — En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.
- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Regresión Lineal Simple (RLS)

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (26.1)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (26.2)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación (??) se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (26.3)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (26.4)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (26.5)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (26.6)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (26.7)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (26.8)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (26.9)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (26.10)$$

Regresión Lineal Simple (RLS)

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (26.11)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (26.12)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación (??) se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (26.13)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (26.14)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (26.15)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (26.16)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (26.17)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (26.18)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (26.19)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (26.20)$$

26.3.3 3. Análisis de Regresión Lineal (RL)

Nota 34. — En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.

3.1 Regresión Lineal Simple (RLS)

- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (26.21)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

3.2 Método de Mínimos Cuadrados

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (26.22)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación ?? se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (26.23)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (26.24)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (26.25)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (26.26)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (26.27)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (26.28)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (26.29)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (26.30)$$

3.3 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Nota 35. – Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo.

– Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right) \\ &= \frac{1}{S_{xx}} E\left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x})\right) \\ &= \frac{1}{S_{xx}} \left[\beta_0 E\left(\sum_{k=1}^n (x_k - \bar{x})\right) + E\left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x})\right) \right. \\ &\quad \left. + E\left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x})\right) \right] = \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1 \end{aligned}$$

por lo tanto

$$E(\hat{\beta}_1) = \beta_1 \quad (26.31)$$

Nota 36. Es decir, $\hat{\beta}_1$ es un estimador insesgado.

Ahora calculemos la varianza:

$$\begin{aligned} V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) \\ &= \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) = \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 \\ &= \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}} \end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad (26.32)$$

Proposición 7.

$$\begin{aligned} E(\hat{\beta}_0) &= \beta_0, \\ V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \\ Cov(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{\sigma^2 \bar{x}}{S_{xx}}. \end{aligned}$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$$e_k = y_k - \hat{y}_k, \text{ se le denomina } \mathbf{residuo}.$$

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \quad (26.33)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$\begin{aligned} SC_E &= \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \\ E(SC_E) &= (n-2)\sigma^2, \text{ por lo tanto} \\ \hat{\sigma}^2 &= \frac{SC_E}{n-2} = MC_E \text{ es un estimador insesgado de } \sigma^2. \end{aligned}$$

3.4 Prueba de Hipótesis en RLS

- Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza.
- Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$.

Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_1, \sigma^2)$.

De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??). Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (26.34)$$

que se distribuye t con $n - 2$ grados de libertad bajo $H_0 : \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (26.35)$$

Para β_0 se puede proceder de manera análoga para

$$H_0 : \beta_0 = \beta_{0,0},$$

$$H_1 : \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right]}, \quad (26.36)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (26.37)$$

- No rechazar $H_0 : \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0 : \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0 : \beta_1 = 0$ puede realizarse de la siguiente manera:

$$S_{yy} = \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2$$

$$\begin{aligned} S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 \\ &= \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\ &= \sum_{k=1}^n \left[(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2 \right] \\ &= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \end{aligned}$$

$$\begin{aligned}
& \sum_{k=1}^n (\hat{y}_k - \bar{y}) (y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) \\
&= \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\
&= \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&= \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&= \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0.
\end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (26.38)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (26.39)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (26.40)$$

Por lo tanto la ecuación (??) se puede reescribir como

$$S_{yy} = SC_R + SC_E \quad (26.41)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$

$$\begin{aligned}
S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}) \\
S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R
\end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 8.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx} \quad (26.42)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E}$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$.

El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n-2$	MC_E	
Total	S_{yy}	$n-1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}} \quad (26.43)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

Estimación de Intervalos en RLS

- Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros.
- El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión.

Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n-2$ grados de libertad. Por tanto un intervalo de confianza de $100(1-\alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \quad (26.44)$$

De igual manera, para β_0 un intervalo de confianza al $100(1-\alpha)\%$ es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (26.45)$$

26.3.4 3. Análisis de Regresión Lineal (RL)

Nota 37. – En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.

3.1 Regresión Lineal Simple (RLS)

- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (26.46)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

3.2 Método de Mínimos Cuadrados

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (26.47)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación ?? se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (26.48)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (26.49)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (26.50)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (26.51)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (26.52)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (26.53)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (26.54)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (26.55)$$

3.3 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Nota 38. – Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo.

– Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right) \\ &= \frac{1}{S_{xx}} E\left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x})\right) \\ &= \frac{1}{S_{xx}} \left[\beta_0 E\left(\sum_{k=1}^n (x_k - \bar{x})\right) + E\left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x})\right) \right. \\ &\quad \left. + E\left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x})\right) \right] = \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1 \end{aligned}$$

por lo tanto

$$E(\hat{\beta}_1) = \beta_1 \quad (26.56)$$

Nota 39. Es decir, $\hat{\beta}_1$ es un estimador insesgado.

Ahora calculemos la varianza:

$$\begin{aligned} V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) \\ &= \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) = \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 \\ &= \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}} \end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad (26.57)$$

Proposición 9.

$$\begin{aligned} E(\hat{\beta}_0) &= \beta_0, \\ V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \\ Cov(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{\sigma^2 \bar{x}}{S_{xx}}. \end{aligned}$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$$e_k = y_k - \hat{y}_k, \text{ se le denomina } \mathbf{residuo}.$$

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \quad (26.58)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$\begin{aligned} SC_E &= \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \\ E(SC_E) &= (n-2)\sigma^2, \text{ por lo tanto} \\ \hat{\sigma}^2 &= \frac{SC_E}{n-2} = MC_E \text{ es un estimador insesgado de } \sigma^2. \end{aligned}$$

3.4 Prueba de Hipótesis en RLS

- Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza.
- Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$.

Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_1, \sigma^2)$.

De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??).

Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (26.59)$$

que se distribuye t con $n - 2$ grados de libertad bajo $H_0: \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (26.60)$$

Para β_0 se puede proceder de manera análoga para

$$H_0: \beta_0 = \beta_{0,0},$$

$$H_1: \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right]}, \quad (26.61)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (26.62)$$

- No rechazar $H_0: \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0: \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0: \beta_1 = 0$ puede realizarse de la siguiente manera:

$$S_{yy} = \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2$$

$$\begin{aligned} S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 \\ &= \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\ &= \sum_{k=1}^n \left[(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2 \right] \\ &= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \end{aligned}$$

$$\begin{aligned}
& \sum_{k=1}^n (\hat{y}_k - \bar{y}) (y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) \\
&= \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\
&= \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&= \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&= \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\
&\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0.
\end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (26.63)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (26.64)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (26.65)$$

Por lo tanto la ecuación (??) se puede reescribir como

$$S_{yy} = SC_R + SC_E \quad (26.66)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$

$$\begin{aligned}
S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}) \\
S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R
\end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 10.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx} \quad (26.67)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E}$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$.

El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n-2$	MC_E	
Total	S_{yy}	$n-1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}} \quad (26.68)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

Estimación de Intervalos en RLS

- Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros.
- El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión.

Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n-2$ grados de libertad. Por tanto un intervalo de confianza de $100(1-\alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}}. \quad (26.69)$$

De igual manera, para β_0 un intervalo de confianza al $100(1-\alpha)\%$ es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (26.70)$$

Predicción

Supongamos que se tiene un valor x_0 de interés, entonces la estimación puntual de este nuevo valor

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 \quad (26.71)$$

Nota 40. Esta nueva observación es independiente de las utilizadas para obtener el modelo de regresión, por tanto, el intervalo en torno a la recta de regresión es inapropiado, puesto que se basa únicamente en los datos empleados para ajustar el modelo de regresión.

El intervalo de confianza en torno a la recta de regresión se refiere a la respuesta media verdadera $x = x_0$, no a observaciones futuras.

Sea y_0 la observación futura en $x = x_0$, y sea \hat{y}_0 dada en la ecuación anterior, el estimador de y_0 . Si se define la variable aleatoria

$$w = y_0 - \hat{y}_0,$$

esta se distribuye normalmente con media cero y varianza

$$V(w) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]$$

dado que y_0 es independiente de \hat{y}_0 , por lo tanto el intervalo de predicción al nivel α para futuras observaciones x_0 es

$$\begin{aligned} \hat{y}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]} &\leq y_0 \\ &\leq \hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]}. \end{aligned}$$

Prueba de falta de ajuste

Es común encontrar que el modelo ajustado no satisface totalmente el modelo necesario para los datos, en este caso es preciso saber qué tan bueno es el modelo propuesto. Para esto se propone la siguiente prueba de hipótesis:

H_0 : El modelo propuesto se ajusta adecuadamente a los datos.

H_1 : El modelo NO se ajusta a los datos.

La prueba implica dividir la suma de cuadrados del error o del residuo en las siguientes dos componentes:

$$SC_E = SC_{EP} + SC_{FDA}$$

donde SC_{EP} es la suma de cuadrados atribuibles al error puro, y SC_{FDA} es la suma de cuadrados atribuible a la falta de ajuste del modelo.

Coeficiente de Determinación

La cantidad

$$R^2 = \frac{SC_R}{S_{yy}} = 1 - \frac{SC_E}{S_{yy}} \quad (26.72)$$

se denomina coeficiente de determinación y se utiliza para saber si el modelo de regresión es suficiente o no. Se puede demostrar que $0 \leq R^2 \leq 1$, una manera de interpretar este valor es que si $R^2 = k$, entonces el modelo de regresión explica el $k * 100\%$ de la variabilidad en los datos. R^2

- No mide la magnitud de la pendiente de la recta de regresión
- Un valor grande de R^2 no implica una pendiente empinada.
- No mide la suficiencia del modelo.
- Valores grandes de R^2 no implican necesariamente que el modelo de regresión proporcionará predicciones precisas para futuras observaciones.

Chapter 27

Regresión Logística

27.1 Conceptos Básicos de la Regresión Logística

La regresión logística es una técnica de modelado estadístico utilizada para predecir la probabilidad de un evento binario (es decir, un evento que tiene dos posibles resultados) en función de una o más variables independientes. A diferencia de la regresión lineal, que se utiliza para predecir valores continuos, la regresión logística se usa cuando la variable dependiente es categórica.

27.2 Diferencias entre Regresión Lineal y Logística

27.2.1 Regresión Lineal

La regresión lineal busca modelar la relación entre una variable dependiente continua Y y una o más variables independientes X_1, X_2, \dots, X_n mediante una ecuación de la forma:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

donde $\beta_0, \beta_1, \dots, \beta_n$ son los coeficientes del modelo y ϵ es el término de error.

27.2.2 Regresión Logística

La regresión logística, en cambio, modela la probabilidad de que un evento ocurra (por ejemplo, éxito vs. fracaso) utilizando la función logística. La variable dependiente Y es binaria, tomando valores de 0 o 1. La ecuación de la regresión logística es:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

donde p es la probabilidad de que $Y = 1$. La función logística es:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$

27.3 Casos de Uso de la Regresión Logística

La regresión logística se utiliza en una variedad de campos para problemas de clasificación binaria, tales como:

- **Medicina:** Predicción de la presencia o ausencia de una enfermedad.
- **Marketing:** Determinación de la probabilidad de que un cliente compre un producto.
- **Finanzas:** Evaluación del riesgo de crédito, es decir, si un cliente va a incumplir o no con un préstamo.
- **Seguridad:** Detección de fraudes o intrusiones.

27.4 Implementación Básica en R

Para implementar una regresión logística en R, primero es necesario instalar y cargar los paquetes necesarios. Aquí se muestra un ejemplo básico de implementación:

27.4.1 Instalación y Configuración de R y RStudio

- Descargue e instale R desde <https://cran.r-project.org/>.
- Descargue e instale RStudio desde <https://rstudio.com/products/rstudio/download/>.

27.4.2 Introducción Básica a R

- Sintaxis básica de R.
- Operaciones básicas: asignación, operaciones aritméticas, funciones básicas.

27.4.3 Ejemplo de Regresión Logística en R

```
# Instalación del paquete necesario
install.packages("stats")

# Carga del paquete
library(stats)

# Ejemplo de conjunto de datos
data <- data.frame(
  outcome = c(1, 0, 1, 0, 1, 1, 0, 1, 0, 0),
  predictor = c(2.3, 1.9, 3.1, 2.8, 3.6, 2.4, 2.1, 3.3, 2.2, 1.7)
)

# Ajuste del modelo de regresión logística
model <- glm(outcome ~ predictor, data = data, family = binomial)

# Resumen del modelo
summary(model)
```

En este ejemplo, se utiliza el conjunto de datos ‘data’ que contiene una variable de resultado binaria ‘outcome’ y una variable predictora continua ‘predictor’. El modelo de regresión logística se ajusta utilizando la función `glm` con la familia binomial.

27.5 Conceptos Básicos

La regresión logística es una técnica de modelado estadístico utilizada para predecir la probabilidad de un evento binario (es decir, un evento que tiene dos posibles resultados) en función de una o más variables independientes. Es ampliamente utilizada en diversas disciplinas, como medicina, economía, biología, y ciencias sociales, para analizar y predecir resultados binarios. Un modelo de regresión logística describe cómo una variable dependiente binaria Y (que puede tomar los valores 0 o 1) está relacionada con una o más variables independientes X_1, X_2, \dots, X_n . A diferencia de la regresión lineal, que predice un valor continuo, la regresión logística predice una probabilidad que puede ser interpretada como la probabilidad de que $Y = 1$ dado un conjunto de valores para X_1, X_2, \dots, X_n .

27.6 Regresión Lineal

La regresión lineal es utilizada para predecir el valor de una variable dependiente continua en función de una o más variables independientes. El modelo de regresión lineal tiene la forma:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon \quad (27.1)$$

donde:

- Y es la variable dependiente.
- β_0 es la intersección con el eje Y o término constante.
- $\beta_1, \beta_2, \dots, \beta_n$ son los coeficientes que representan la relación entre las variables independientes y la variable dependiente.
- X_1, X_2, \dots, X_n son las variables independientes.
- ϵ es el término de error, que representa la desviación de los datos observados de los valores predichos por el modelo.

El objetivo de la regresión lineal es encontrar los valores de los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ que minimicen la suma de los cuadrados de las diferencias entre los valores observados y los valores predichos. Este método se conoce como mínimos cuadrados ordinarios (OLS, por sus siglas en inglés). La función de costo a minimizar es:

$$J(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (27.2)$$

donde:

- y_i es el valor observado de la variable dependiente para la i -ésima observación.
- \hat{y}_i es el valor predicho por el modelo para la i -ésima observación, dado por:

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_n x_{in} \quad (27.3)$$

Para encontrar los valores óptimos de los coeficientes, se toman las derivadas parciales de la función de costo con respecto a cada coeficiente y se igualan a cero:

$$\frac{\partial J}{\partial \beta_j} = 0 \quad \text{para } j = 0, 1, \dots, n \quad (27.4)$$

Resolviendo este sistema de ecuaciones, se obtienen los valores de los coeficientes que minimizan la función de costo.

27.7 Regresión Logística

La deducción de la fórmula de la regresión logística comienza con la necesidad de modelar la probabilidad de un evento binario. Queremos encontrar una función que relacione las variables independientes con la probabilidad de que la variable dependiente tome el valor 1. La probabilidad de que el evento ocurra, $P(Y = 1)$, se denota como p . La probabilidad de que el evento no ocurra, $P(Y = 0)$, es $1 - p$. Los *odds* (chances) de que ocurra el evento se definen como:

$$\text{odds} = \frac{p}{1 - p} \quad (27.5)$$

Los *odds* nos indican cuántas veces más probable es que ocurra el evento frente a que no ocurra. Para simplificar el modelado de los *odds*, aplicamos el logaritmo natural, obteniendo la función logit:

$$\text{logit}(p) = \log\left(\frac{p}{1 - p}\right) \quad (27.6)$$

La transformación logit es útil porque convierte el rango de la probabilidad $(0, 1)$ al rango de números reales $(-\infty, \infty)$. La idea clave de la regresión logística es modelar la transformación logit de la probabilidad como una combinación lineal de las variables independientes:

$$\text{logit}(p) = \log\left(\frac{p}{1 - p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad (27.7)$$

Aquí, β_0 es el término constante y $\beta_1, \beta_2, \dots, \beta_n$ son los coeficientes asociados con las variables independientes X_1, X_2, \dots, X_n . Para expresar p en función de una combinación lineal de las variables independientes, invertimos la transformación logit. Partimos de la ecuación:

$$\log\left(\frac{p}{1 - p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

Aplicamos la exponenciación a ambos lados:

$$\frac{p}{1 - p} = e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n}$$

Despejamos p :

$$p = \frac{e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n}}{1 + e^{\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n}}$$

La expresión final que obtenemos es conocida como la función logística:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}} \quad (27.8)$$

Esta función describe cómo las variables independientes se relacionan con la probabilidad de que el evento de interés ocurra. Los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ se estiman a partir de los datos utilizando el método de máxima verosimilitud.

27.8 Método de Máxima Verosimilitud

Para estimar los coeficientes $\beta_0, \beta_1, \dots, \beta_n$ en la regresión logística, utilizamos el método de máxima verosimilitud. La idea es encontrar los valores de los coeficientes que maximicen la probabilidad de observar los datos dados. Esta probabilidad se expresa mediante la función de verosimilitud L . La función de verosimilitud $L(\beta_0, \beta_1, \dots, \beta_n)$ para un conjunto de n observaciones se define como el producto de las probabilidades de las observaciones dadas las variables independientes:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (27.9)$$

donde:

- p_i es la probabilidad predicha de que $Y_i = 1$,
- y_i es el valor observado de la variable dependiente para la i -ésima observación.

Trabajar directamente con esta función de verosimilitud puede ser complicado debido al producto de muchas probabilidades, especialmente si n es grande. Para simplificar los cálculos, se utiliza el logaritmo de la función de verosimilitud, conocido como la función de log-verosimilitud. El uso del logaritmo simplifica significativamente la diferenciación y maximización de la función. La función de log-verosimilitud se define como:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (27.10)$$

Aquí, \log representa el logaritmo natural. Esta transformación es válida porque el logaritmo es una función monótona creciente, lo que significa que maximizar la log-verosimilitud es equivalente a maximizar la verosimilitud original. En la regresión logística, la probabilidad p_i está dada por la función logística:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \quad (27.11)$$

Sustituyendo esta expresión en la función de log-verosimilitud, obtenemos:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n \left[y_i \log \left(\frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \right) + \right. \\ &\quad \left. (1 - y_i) \log \left(1 - \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \right) \right] \end{aligned}$$

Simplificando esta expresión, notamos que:

$$\log \left(\frac{1}{1 + e^{-z}} \right) = -\log(1 + e^{-z})$$

y

$$\log \left(1 - \frac{1}{1 + e^{-z}} \right) = \log \left(\frac{e^{-z}}{1 + e^{-z}} \right) = -z - \log(1 + e^{-z})$$

Aplicando estas identidades, la función de log-verosimilitud se convierte en:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n \left[y_i (-\log(1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})})) + \right. \\ &\quad \left. (1 - y_i) \left(-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}) \right) \right] \end{aligned}$$

Simplificando aún más, obtenemos:

$$\begin{aligned} \log L(\beta_0, \beta_1, \dots, \beta_n) &= \sum_{i=1}^n [y_i(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) \\ &\quad - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}})] \end{aligned}$$

Para simplificar aún más la notación, podemos utilizar notación matricial. Definimos la matriz \mathbf{X} de tamaño $n \times (k + 1)$ y el vector de coeficientes $\boldsymbol{\beta}$ de tamaño $(k + 1) \times 1$ como sigue:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \quad (27.12)$$

Entonces, la expresión para la función de log-verosimilitud es:

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n \left[y_i (\mathbf{X}_i \boldsymbol{\beta}) - \log(1 + e^{\mathbf{X}_i \boldsymbol{\beta}}) \right] \quad (27.13)$$

donde \mathbf{X}_i es la i -ésima fila de la matriz \mathbf{X} . Esta notación matricial simplifica la implementación y la derivación de los estimadores de los coeficientes en la regresión logística. Utilizando métodos numéricos, como el algoritmo de Newton-Raphson, se pueden encontrar los coeficientes que maximizan la función de log-verosimilitud. Para maximizar la función de log-verosimilitud, derivamos esta función con respecto a cada uno de los coeficientes β_j y encontramos los puntos críticos. La derivada parcial de la función de log-verosimilitud con respecto a β_j es:

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = \sum_{i=1}^n \left[y_i X_{ij} - \frac{X_{ij} e^{\mathbf{X}_i \boldsymbol{\beta}}}{1 + e^{\mathbf{X}_i \boldsymbol{\beta}}} \right] \quad (27.14)$$

Simplificando, esta derivada se puede expresar como:

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = \sum_{i=1}^n X_{ij} (y_i - p_i), \quad \text{donde } p_i = \frac{1}{1 + e^{-\mathbf{X}_i \boldsymbol{\beta}}} \quad (27.15)$$

Para encontrar los coeficientes que maximizan la log-verosimilitud, resolvemos el sistema de ecuaciones

$$\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j} = 0 \quad \text{para todos los } j = 0, 1, \dots, k.$$

Este sistema de ecuaciones no tiene una solución analítica cerrada, por lo que se resuelve numéricamente utilizando métodos iterativos como el algoritmo de Newton-Raphson.

27.9 Método de Newton-Raphson

El método de Newton-Raphson es un algoritmo iterativo que se utiliza para encontrar las raíces de una función. En el contexto de la regresión logística, se utiliza para maximizar la función de log-verosimilitud encontrando los valores de los coeficientes $\beta_0, \beta_1, \dots, \beta_n$. Este método se basa en una aproximación de segundo orden de la función objetivo. Dado un valor inicial de los coeficientes $\beta^{(0)}$, se actualiza iterativamente el valor de los coeficientes utilizando la fórmula:

$$\beta^{(t+1)} = \beta^{(t)} - [\mathbf{H}(\beta^{(t)})]^{-1} \nabla \log L(\beta^{(t)}) \quad (27.16)$$

donde:

- $\beta^{(t)}$ es el vector de coeficientes en la t -ésima iteración.
- $\nabla \log L(\beta^{(t)})$ es el gradiente de la función de log-verosimilitud con respecto a los coeficientes β :

$$\nabla \log L(\beta) = \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \quad (27.17)$$

donde \mathbf{y} es el vector de valores observados y \mathbf{p} es el vector de probabilidades.

- $\mathbf{H}(\beta^{(t)})$ es la matriz Hessiana (matriz de segundas derivadas) evaluada en $\beta^{(t)}$:

$$\mathbf{H}(\beta) = -\mathbf{X}^T \mathbf{W} \mathbf{X} \quad (27.18)$$

donde \mathbf{W} es una matriz diagonal de pesos con elementos $w_i = p_i(1 - p_i)$.

En resumen:

Algoritmo 5. *El algoritmo Newton-Raphson para la regresión logística se puede resumir en los siguientes pasos:*

1. Inicializar el vector de coeficientes $\beta^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).
2. Calcular el gradiente $\nabla \log L(\beta^{(t)})$ y la matriz Hessiana $\mathbf{H}(\beta^{(t)})$ en la iteración t .
3. Actualizar los coeficientes utilizando la fórmula:

$$\beta^{(t+1)} = \beta^{(t)} - [\mathbf{H}(\beta^{(t)})]^{-1} \nabla \log L(\beta^{(t)}) \quad (27.19)$$

4. Repetir los pasos 2 y 3 hasta que la diferencia entre $\beta^{(t+1)}$ y $\beta^{(t)}$ sea menor que un umbral pre-definido (criterio de convergencia).

En resumen, el método de Newton-Raphson permite encontrar los coeficientes que maximizan la función de log-verosimilitud de manera eficiente.

27.10 Especificando

En específico para un conjunto de n observaciones, la función de verosimilitud L se define como el producto de las probabilidades individuales de observar cada dato:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (27.20)$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación y p_i es la probabilidad predicha de que $Y_i = 1$. Aquí, p_i es dado por la función logística:

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in})}} \quad (27.21)$$

Tomando el logaritmo:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (27.22)$$

Sustituyendo p_i :

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}) \right] \quad (27.23)$$

Dado que el objetivo es encontrar los valores de $\beta_0, \beta_1, \dots, \beta_n$ que maximicen la función de log-verosimilitud. Para β_j , la derivada parcial de la función de log-verosimilitud es:

$$\frac{\partial \log L}{\partial \beta_j} = \sum_{i=1}^n \left[y_i X_{ij} - \frac{X_{ij} e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}}{1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}} \right] \quad (27.24)$$

Esto se simplifica a (comparar con la ecuación ??):

$$\frac{\partial \log L}{\partial \beta_j} = \sum_{i=1}^n X_{ij} (y_i - p_i) \quad (27.25)$$

Para maximizar la log-verosimilitud, resolvemos el sistema de ecuaciones $\frac{\partial \log L}{\partial \beta_j} = 0$ para todos los j de 0 a n , mismo que se resuelve numéricamente utilizando métodos el algoritmo de Newton-Raphson. El método de Newton-Raphson se basa en una aproximación de segundo orden de la función objetivo. Dado un valor inicial de los coeficientes $\beta^{(0)}$, se iterativamente actualiza el valor de los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - [\mathbf{H}(\beta^{(k)})]^{-1} \mathbf{g}(\beta^{(k)}) \quad (27.26)$$

donde:

- $\beta^{(k)}$ es el vector de coeficientes en la k -ésima iteración.
- $\mathbf{g}(\beta^{(k)})$ es el gradiente (vector de primeras derivadas) evaluado en $\beta^{(k)}$:

$$\mathbf{g}(\beta) = \frac{\partial \log L}{\partial \beta} = \sum_{i=1}^n \mathbf{X}_i (y_i - p_i) \quad (27.27)$$

donde \mathbf{X}_i es el vector de valores de las variables independientes para la i -ésima observación (comparar con ecuación ??).

- $\mathbf{H}(\beta^{(k)})$ es la matriz Hessiana (matriz de segundas derivadas) evaluada en $\beta^{(k)}$:

$$\mathbf{H}(\beta) = \frac{\partial^2 \log L}{\partial \beta \partial \beta^T} = - \sum_{i=1}^n p_i (1 - p_i) \mathbf{X}_i \mathbf{X}_i^T, \quad (27.28)$$

comparar con ecuación ??

Algoritmo 6. Los pasos del algoritmo Newton-Raphson para la regresión logística son:

1. Inicializar el vector de coeficientes $\beta^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).
2. Calcular el gradiente $\mathbf{g}(\beta^{(k)})$ y la matriz Hessiana $\mathbf{H}(\beta^{(k)})$ en la iteración k .
3. Actualizar los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - \left[\mathbf{H}(\beta^{(k)}) \right]^{-1} \mathbf{g}(\beta^{(k)}) \quad (27.29)$$

4. Repetir los pasos 2 y 3 hasta que la diferencia entre $\beta^{(k+1)}$ y $\beta^{(k)}$ sea menor que un umbral pre-definido (criterio de convergencia).

Como se puede observar la diferencia entre el Algoritmo ?? y el Algoritmo ?? son mínimas

Notas finales

En el contexto de la regresión logística, los vectores X_1, X_2, \dots, X_n representan las variables independientes. Cada X_j es un vector columna que contiene los valores de la variable independiente j para cada una de las n observaciones. Es decir,

$$X_j = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{bmatrix} \quad (27.30)$$

Para simplificar la notación y los cálculos, a menudo combinamos todos los vectores de variables independientes en una única matriz de diseño \mathbf{X} de tamaño $n \times (k+1)$, donde n es el número de observaciones y $k+1$ es el número de variables independientes más el término de intercepto. La primera columna de \mathbf{X} corresponde a un vector de unos para el término de intercepto, y las demás columnas corresponden a los valores de las variables independientes:

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix} \quad (27.31)$$

revisar la ecuación ?? . De esta forma, el modelo logit puede ser escrito de manera compacta utilizando la notación matricial:

$$\text{logit}(p) = \log \left(\frac{p}{1-p} \right) = \mathbf{X}\beta \quad (27.32)$$

donde β es el vector de coeficientes:

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \quad (27.33)$$

Así, la probabilidad p se puede expresar como:

$$p = \frac{1}{1 + e^{-\mathbf{X}\boldsymbol{\beta}}} \quad (27.34)$$

Comparar la ecuación anterior con la ecuación ???. Esta notación matricial simplifica la implementación y la derivación de los estimadores de los coeficientes en la regresión logística. Para estimar los coeficientes $\boldsymbol{\beta}$ en la regresión logística, se utiliza el método de máxima verosimilitud. La función de verosimilitud $L(\boldsymbol{\beta})$ se define como el producto de las probabilidades de las observaciones dadas las variables independientes, recordemos la ecuación ???:

$$L(\boldsymbol{\beta}) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i} \quad (27.35)$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación, y p_i es la probabilidad predicha de que $Y_i = 1$. La función de log-verosimilitud, que es más fácil de maximizar, se obtiene tomando el logaritmo natural de la función de verosimilitud (??):

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)] \quad (27.36)$$

Sustituyendo $p_i = \frac{1}{1+e^{-\mathbf{X}_i\boldsymbol{\beta}}}$, donde \mathbf{X}_i es la i -ésima fila de la matriz de diseño \mathbf{X} , obtenemos:

$$\log L(\boldsymbol{\beta}) = \sum_{i=1}^n \left[y_i (\mathbf{X}_i \boldsymbol{\beta}) - \log(1 + e^{\mathbf{X}_i \boldsymbol{\beta}}) \right] \quad (27.37)$$

Para encontrar los valores de $\boldsymbol{\beta}$ que maximizan la función de log-verosimilitud, se utiliza un algoritmo iterativo como el método de Newton-Raphson. Este método requiere calcular el gradiente y la matriz Hessiana de la función de log-verosimilitud.

El gradiente de la función de log-verosimilitud con respecto a $\boldsymbol{\beta}$ es (?? y ??):

$$\nabla \log L(\boldsymbol{\beta}) = \mathbf{X}^T (\mathbf{y} - \mathbf{p}) \quad (27.38)$$

donde \mathbf{y} es el vector de valores observados y \mathbf{p} es el vector de probabilidades predichas.

La matriz Hessiana de la función de log-verosimilitud es (?? y ??):

$$\mathbf{H}(\boldsymbol{\beta}) = -\mathbf{X}^T \mathbf{W} \mathbf{X} \quad (27.39)$$

donde \mathbf{W} es una matriz diagonal de pesos con elementos $w_i = p_i(1 - p_i)$.

El método de Newton-Raphson actualiza los coeficientes $\boldsymbol{\beta}$ de la siguiente manera:

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - [\mathbf{H}(\boldsymbol{\beta}^{(t)})]^{-1} \nabla \log L(\boldsymbol{\beta}^{(t)}) \quad (27.40)$$

Iterando este proceso hasta que la diferencia entre $\boldsymbol{\beta}^{(t+1)}$ y $\boldsymbol{\beta}^{(t)}$ sea menor que un umbral predefinido (??, ??, ?? y ??), se obtienen los estimadores de máxima verosimilitud para los coeficientes de la regresión logística.

Chapter 28

Elementos de Probabilidad

28.1 Introducción

Los fundamentos de probabilidad y estadística son esenciales para comprender y aplicar técnicas de análisis de datos y modelado estadístico, incluyendo la regresión lineal y logística. Este capítulo proporciona una revisión de los conceptos clave en probabilidad y estadística que son relevantes para estos métodos.

28.2 Probabilidad

La probabilidad es una medida de la incertidumbre o el grado de creencia en la ocurrencia de un evento. Los conceptos fundamentales incluyen:

28.2.1 Espacio Muestral y Eventos

El espacio muestral, denotado como S , es el conjunto de todos los posibles resultados de un experimento aleatorio. Un evento es un subconjunto del espacio muestral. Por ejemplo, si lanzamos un dado, el espacio muestral es:

$$S = \{1, 2, 3, 4, 5, 6\}$$

Un evento podría ser obtener un número par:

$$E = \{2, 4, 6\}$$

28.2.2 Definiciones de Probabilidad

Existen varias definiciones de probabilidad, incluyendo la probabilidad clásica, la probabilidad frecuentista y la probabilidad bayesiana.

Probabilidad Clásica

La probabilidad clásica se define como el número de resultados favorables dividido por el número total de resultados posibles:

$$P(E) = \frac{|E|}{|S|}$$

donde $|E|$ es el número de elementos en el evento E y $|S|$ es el número de elementos en el espacio muestral S .

Probabilidad Frecuentista

La probabilidad frecuentista se basa en la frecuencia relativa de ocurrencia de un evento en un gran número de repeticiones del experimento:

$$P(E) = \lim_{n \rightarrow \infty} \frac{n_E}{n}$$

donde n_E es el número de veces que ocurre el evento E y n es el número total de repeticiones del experimento.

Probabilidad Bayesiana

La probabilidad bayesiana se interpreta como un grado de creencia actualizado a medida que se dispone de nueva información. Se basa en el teorema de Bayes:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

donde $P(A|B)$ es la probabilidad de A dado B , $P(B|A)$ es la probabilidad de B dado A , $P(A)$ y $P(B)$ son las probabilidades de A y B respectivamente.

28.3 Estadística Bayesiana

La estadística bayesiana proporciona un enfoque coherente para el análisis de datos basado en el teorema de Bayes. Los conceptos fundamentales incluyen:

28.3.1 Prior y Posterior

Distribución Prior

La distribución prior (apriori) representa nuestra creencia sobre los parámetros antes de observar los datos. Es una distribución de probabilidad que refleja nuestra incertidumbre inicial sobre los parámetros. Por ejemplo, si creemos que un parámetro θ sigue una distribución normal con media μ_0 y varianza σ_0^2 , nuestra prior sería:

$$P(\theta) = \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(\theta-\mu_0)^2}{2\sigma_0^2}}$$

Verosimilitud

La verosimilitud (likelihood) es la probabilidad de observar los datos dados los parámetros. Es una función de los parámetros θ dada una muestra de datos X :

$$L(\theta; X) = P(X|\theta)$$

donde X son los datos observados y θ son los parámetros del modelo.

Distribución Posterior

La distribución posterior (a posteriori) combina la información de la prior y la verosimilitud utilizando el teorema de Bayes. Representa nuestra creencia sobre los parámetros después de observar los datos:

$$P(\theta|X) = \frac{P(X|\theta)P(\theta)}{P(X)}$$

donde $P(\theta|X)$ es la distribución posterior, $P(X|\theta)$ es la verosimilitud, $P(\theta)$ es la prior y $P(X)$ es la probabilidad marginal de los datos.

La probabilidad marginal de los datos $P(X)$ se puede calcular como:

$$P(X) = \int_{\Theta} P(X|\theta)P(\theta)d\theta$$

donde Θ es el espacio de todos los posibles valores del parámetro θ .

28.4 Distribuciones de Probabilidad

Las distribuciones de probabilidad describen cómo se distribuyen los valores de una variable aleatoria. Existen distribuciones de probabilidad discretas y continuas.

28.4.1 Distribuciones Discretas

Una variable aleatoria discreta toma un número finito o contable de valores. Algunas distribuciones discretas comunes incluyen:

Distribución Binomial

La distribución binomial describe el número de éxitos en una serie de ensayos de Bernoulli independientes y con la misma probabilidad de éxito. La función de probabilidad es:

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

donde X es el número de éxitos, n es el número de ensayos, p es la probabilidad de éxito en cada ensayo, y $\binom{n}{k}$ es el coeficiente binomial.

La función generadora de momentos (MGF) para la distribución binomial es:

$$M_X(t) = (1 - p + pe^t)^n$$

El valor esperado y la varianza de una variable aleatoria binomial son:

$$\begin{aligned} E(X) &= np \\ \text{Var}(X) &= np(1 - p) \end{aligned}$$

Distribución de Poisson

La distribución de Poisson describe el número de eventos que ocurren en un intervalo de tiempo fijo o en un área fija. La función de probabilidad es:

$$P(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}$$

donde X es el número de eventos, λ es la tasa media de eventos por intervalo, y k es el número de eventos observados.

La función generadora de momentos (MGF) para la distribución de Poisson es:

$$M_X(t) = e^{\lambda(e^t - 1)}$$

El valor esperado y la varianza de una variable aleatoria de Poisson son:

$$\begin{aligned} E(X) &= \lambda \\ \text{Var}(X) &= \lambda \end{aligned}$$

28.4.2 Distribuciones Continuas

Una variable aleatoria continua toma un número infinito de valores en un intervalo continuo. Algunas distribuciones continuas comunes incluyen:

Distribución Normal

La distribución normal, también conocida como distribución gaussiana, es una de las distribuciones más importantes en estadística. La función de densidad de probabilidad es:

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

donde x es un valor de la variable aleatoria, μ es la media, y σ es la desviación estándar.

La función generadora de momentos (MGF) para la distribución normal es:

$$M_X(t) = e^{\mu t + \frac{1}{2}\sigma^2 t^2}$$

El valor esperado y la varianza de una variable aleatoria normal son:

$$\begin{aligned} E(X) &= \mu \\ \text{Var}(X) &= \sigma^2 \end{aligned}$$

Distribución Exponencial

La distribución exponencial describe el tiempo entre eventos en un proceso de Poisson. La función de densidad de probabilidad es:

$$f(x) = \lambda e^{-\lambda x}$$

donde x es el tiempo entre eventos y λ es la tasa media de eventos.

La función generadora de momentos (MGF) para la distribución exponencial es:

$$M_X(t) = \frac{\lambda}{\lambda - t}, \quad \text{para } t < \lambda$$

El valor esperado y la varianza de una variable aleatoria exponencial son:

$$\begin{aligned} E(X) &= \frac{1}{\lambda} \\ \text{Var}(X) &= \frac{1}{\lambda^2} \end{aligned}$$

28.5 Estadística Descriptiva

La estadística descriptiva resume y describe las características de un conjunto de datos. Incluye medidas de tendencia central, medidas de dispersión y medidas de forma.

28.5.1 Medidas de Tendencia Central

Las medidas de tendencia central incluyen la media, la mediana y la moda.

Media

La media aritmética es la suma de los valores dividida por el número de valores:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

donde x_i son los valores de la muestra y n es el tamaño de la muestra.

Mediana

La mediana es el valor medio cuando los datos están ordenados. Si el número de valores es impar, la mediana es el valor central. Si es par, es el promedio de los dos valores centrales.

Moda

La moda es el valor que ocurre con mayor frecuencia en un conjunto de datos.

28.5.2 Medidas de Dispersión

Las medidas de dispersión incluyen el rango, la varianza y la desviación estándar.

Rango

El rango es la diferencia entre el valor máximo y el valor mínimo de los datos:

$$\text{Rango} = x_{\max} - x_{\min}$$

Varianza

La varianza es la media de los cuadrados de las diferencias entre los valores y la media:

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

Desviación Estándar

La desviación estándar es la raíz cuadrada de la varianza:

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

28.6 Inferencia Estadística

La inferencia estadística es el proceso de sacar conclusiones sobre una población a partir de una muestra. Incluye la estimación de parámetros y la prueba de hipótesis.

28.6.1 Estimación de Parámetros

La estimación de parámetros implica el uso de datos muestrales para estimar los parámetros de una población.

Estimador Puntual

Un estimador puntual proporciona un único valor como estimación de un parámetro de la población. Por ejemplo, la media muestral \bar{x} es un estimador puntual de la media poblacional μ . Otros ejemplos de estimadores puntuales son:

- **Mediana muestral** (\tilde{x}): Estimador de la mediana poblacional.
- **Varianza muestral** (s^2): Estimador de la varianza poblacional σ^2 , definido como:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

- **Desviación estándar muestral** (s): Estimador de la desviación estándar poblacional σ , definido como:

$$s = \sqrt{s^2}$$

Propiedades de los Estimadores Puntuales

Los estimadores puntuales deben cumplir ciertas propiedades deseables, como:

- **Insesgadez**: Un estimador es insesgado si su valor esperado es igual al valor del parámetro que estima.

$$E(\hat{\theta}) = \theta$$

- **Consistencia:** Un estimador es consistente si converge en probabilidad al valor del parámetro a medida que el tamaño de la muestra tiende a infinito.
- **Eficiencia:** Un estimador es eficiente si tiene la varianza más baja entre todos los estimadores insesgados.

Estimador por Intervalo

Un estimador por intervalo proporciona un rango de valores dentro del cual se espera que se encuentre el parámetro poblacional con un cierto nivel de confianza. Por ejemplo, un intervalo de confianza para la media es:

$$\left(\bar{x} - z \frac{\sigma}{\sqrt{n}}, \bar{x} + z \frac{\sigma}{\sqrt{n}} \right)$$

donde z es el valor crítico correspondiente al nivel de confianza deseado, σ es la desviación estándar poblacional y n es el tamaño de la muestra.

28.6.2 Prueba de Hipótesis

La prueba de hipótesis es un procedimiento para decidir si una afirmación sobre un parámetro poblacional es consistente con los datos muestrales.

Hipótesis Nula y Alternativa

La hipótesis nula (H_0) es la afirmación que se somete a prueba, y la hipótesis alternativa (H_a) es la afirmación que se acepta si se rechaza la hipótesis nula.

Nivel de Significancia

El nivel de significancia (α) es la probabilidad de rechazar la hipótesis nula cuando es verdadera. Un valor comúnmente utilizado es $\alpha = 0.05$.

Estadístico de Prueba

El estadístico de prueba es una medida calculada a partir de los datos muestrales que se utiliza para decidir si se rechaza la hipótesis nula. Por ejemplo, en una prueba t para la media:

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

donde \bar{x} es la media muestral, μ_0 es la media poblacional bajo la hipótesis nula, s es la desviación estándar muestral y n es el tamaño de la muestra.

P-valor

El p-valor es la probabilidad de obtener un valor del estadístico de prueba al menos tan extremo como el observado, bajo la suposición de que la hipótesis nula es verdadera. Si el p-valor es menor que el nivel de significancia α , se rechaza la hipótesis nula. El p-valor se interpreta de la siguiente manera:

- **P-valor bajo (p < 0.05):** Evidencia suficiente para rechazar la hipótesis nula.
- **P-valor alto (p > 0.05):** No hay suficiente evidencia para rechazar la hipótesis nula.

Tipos de Errores

En la prueba de hipótesis, se pueden cometer dos tipos de errores:

- **Error Tipo I** (α): Rechazar la hipótesis nula cuando es verdadera.
- **Error Tipo II** (β): No rechazar la hipótesis nula cuando es falsa.

Tabla de Errores en la Prueba de Hipótesis

A continuación se presenta una tabla que muestra los posibles resultados en una prueba de hipótesis, incluyendo los falsos positivos (error tipo I) y los falsos negativos (error tipo II):

	Hipótesis Nula Verdadera	Hipótesis Nula Falsa
Rechazar H_0	Error Tipo I (α)	Aceptar H_a
No Rechazar H_0	Aceptar H_0	Error Tipo II (β)

Table 28.1: Resultados de la Prueba de Hipótesis

Chapter 29

Matemáticas Detrás de la Regresión Logística

29.1 Introducción

La regresión logística es una técnica de modelado estadístico utilizada para predecir la probabilidad de un evento binario en función de una o más variables independientes. Este capítulo profundiza en las matemáticas subyacentes a la regresión logística, incluyendo la función logística, la función de verosimilitud, y los métodos para estimar los coeficientes del modelo.

29.2 Función Logística

La función logística es la base de la regresión logística. Esta función transforma una combinación lineal de variables independientes en una probabilidad.

29.2.1 Definición

La función logística se define como:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$

donde p es la probabilidad de que el evento ocurra, $\beta_0, \beta_1, \dots, \beta_n$ son los coeficientes del modelo, y X_1, X_2, \dots, X_n son las variables independientes.

29.2.2 Propiedades

La función logística tiene varias propiedades importantes:

- **Rango:** La función logística siempre produce un valor entre 0 y 1, lo que la hace adecuada para modelar probabilidades.
- **Monotonía:** La función es monótona creciente, lo que significa que a medida que la combinación lineal de variables independientes aumenta, la probabilidad también aumenta.
- **Simetría:** La función logística es simétrica en torno a $p = 0.5$.

29.3 Función de Verosimilitud

La función de verosimilitud se utiliza para estimar los coeficientes del modelo de regresión logística. Esta función mide la probabilidad de observar los datos dados los coeficientes del modelo.

29.3.1 Definición

Para un conjunto de n observaciones, la función de verosimilitud L se define como el producto de las probabilidades individuales de observar cada dato:

$$L(\beta_0, \beta_1, \dots, \beta_n) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$$

donde y_i es el valor observado de la variable dependiente para la i -ésima observación y p_i es la probabilidad predicha de que $Y_i = 1$.

29.3.2 Función de Log-Verosimilitud

Para simplificar los cálculos, trabajamos con el logaritmo de la función de verosimilitud, conocido como la función de log-verosimilitud. Tomar el logaritmo convierte el producto en una suma:

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

Sustituyendo p_i :

$$\log L(\beta_0, \beta_1, \dots, \beta_n) = \sum_{i=1}^n \left[y_i (\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}) - \log(1 + e^{\beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_n X_{in}}) \right]$$

29.4 Estimación de Coeficientes

Los coeficientes del modelo de regresión logística se estiman maximizando la función de log-verosimilitud. Este proceso generalmente se realiza mediante métodos iterativos como el algoritmo de Newton-Raphson.

29.4.1 Gradiente y Hessiana

Para maximizar la función de log-verosimilitud, necesitamos calcular su gradiente y su matriz Hessiana.

Gradiente

El gradiente de la función de log-verosimilitud con respecto a los coeficientes β es:

$$\mathbf{g}(\beta) = \frac{\partial \log L}{\partial \beta} = \sum_{i=1}^n \mathbf{X}_i (y_i - p_i)$$

donde \mathbf{X}_i es el vector de valores de las variables independientes para la i -ésima observación.

Hessiana

La matriz Hessiana de la función de log-verosimilitud con respecto a los coeficientes β es:

$$\mathbf{H}(\beta) = \frac{\partial^2 \log L}{\partial \beta \partial \beta^T} = - \sum_{i=1}^n p_i(1 - p_i) \mathbf{X}_i \mathbf{X}_i^T$$

29.4.2 Algoritmo Newton-Raphson

El algoritmo Newton-Raphson se utiliza para encontrar los valores de los coeficientes que maximizan la función de log-verosimilitud. El algoritmo se puede resumir en los siguientes pasos:

1. Inicializar el vector de coeficientes $\beta^{(0)}$ (por ejemplo, con ceros o valores pequeños aleatorios).
2. Calcular el gradiente $\mathbf{g}(\beta^{(k)})$ y la matriz Hessiana $\mathbf{H}(\beta^{(k)})$ en la iteración k .
3. Actualizar los coeficientes utilizando la fórmula:

$$\beta^{(k+1)} = \beta^{(k)} - \left[\mathbf{H}(\beta^{(k)}) \right]^{-1} \mathbf{g}(\beta^{(k)})$$

4. Repetir los pasos 2 y 3 hasta que la diferencia entre $\beta^{(k+1)}$ y $\beta^{(k)}$ sea menor que un umbral predefinido (criterio de convergencia).

29.5 Validación del Modelo

Una vez que se han estimado los coeficientes del modelo de regresión logística, es importante validar el modelo para asegurarse de que proporciona predicciones precisas.

29.5.1 Curva ROC y AUC

La curva ROC (Receiver Operating Characteristic) es una herramienta gráfica utilizada para evaluar el rendimiento de un modelo de clasificación binaria. El área bajo la curva (AUC) mide la capacidad del modelo para distinguir entre las clases.

29.5.2 Matriz de Confusión

La matriz de confusión es una tabla que resume el rendimiento de un modelo de clasificación al comparar las predicciones del modelo con los valores reales. Los términos en la matriz de confusión incluyen verdaderos positivos, falsos positivos, verdaderos negativos y falsos negativos.

Chapter 30

Preparación de Datos y Selección de Variables

30.1 Introducción

La preparación de datos y la selección de variables son pasos cruciales en el proceso de modelado estadístico. Un modelo bien preparado y con las variables adecuadas puede mejorar significativamente la precisión y la interpretabilidad del modelo. Este capítulo proporciona una revisión detallada de las técnicas de limpieza de datos, tratamiento de datos faltantes, codificación de variables categóricas y selección de variables.

30.2 Importancia de la Preparación de Datos

La calidad de los datos es fundamental para el éxito de cualquier análisis estadístico. Los datos sin limpiar pueden llevar a modelos inexactos y conclusiones erróneas. La preparación de datos incluye varias etapas:

- Limpieza de datos
- Tratamiento de datos faltantes
- Codificación de variables categóricas
- Selección y transformación de variables

30.3 Limpieza de Datos

La limpieza de datos es el proceso de detectar y corregir (o eliminar) los datos incorrectos, incompletos o irrelevantes. Este proceso incluye:

- Eliminación de duplicados
- Corrección de errores tipográficos
- Consistencia de formato
- Tratamiento de valores extremos (outliers)

30.4 Tratamiento de Datos Faltantes

Los datos faltantes son un problema común en los conjuntos de datos y pueden afectar la calidad de los modelos. Hay varias estrategias para manejar los datos faltantes:

- **Eliminación de Datos Faltantes:** Se eliminan las filas o columnas con datos faltantes.
- **Imputación:** Se reemplazan los valores faltantes con estimaciones, como la media, la mediana o la moda.
- **Modelos Predictivos:** Se utilizan modelos predictivos para estimar los valores faltantes.

30.4.1 Imputación de la Media

Una técnica común es reemplazar los valores faltantes con la media de la variable. Esto se puede hacer de la siguiente manera:

$$x_i = \begin{cases} x_i & \text{si } x_i \text{ no es faltante} \\ \bar{x} & \text{si } x_i \text{ es faltante} \end{cases}$$

donde \bar{x} es la media de la variable.

30.5 Codificación de Variables Categóricas

Las variables categóricas deben ser convertidas a un formato numérico antes de ser usadas en un modelo de regresión logística. Hay varias técnicas para codificar variables categóricas:

30.5.1 Codificación One-Hot

La codificación one-hot crea una columna binaria para cada categoría. Por ejemplo, si tenemos una variable categórica con tres categorías (A, B, C), se crean tres columnas:

$$\begin{aligned} A &= [1, 0, 0] \\ B &= [0, 1, 0] \\ C &= [0, 0, 1] \end{aligned}$$

30.5.2 Codificación Ordinal

La codificación ordinal asigna un valor entero único a cada categoría, preservando el orden natural de las categorías. Por ejemplo:

$$\begin{aligned} \text{Bajo} &= 1 \\ \text{Medio} &= 2 \\ \text{Alto} &= 3 \end{aligned}$$

30.6 Selección de Variables

La selección de variables es el proceso de elegir las variables más relevantes para el modelo. Existen varias técnicas para la selección de variables:

30.6.1 Métodos de Filtrado

Los métodos de filtrado seleccionan variables basadas en criterios estadísticos, como la correlación o la chi-cuadrado. Algunas técnicas comunes incluyen:

- **Análisis de Correlación:** Se seleccionan variables con alta correlación con la variable dependiente y baja correlación entre ellas.
- **Pruebas de Chi-cuadrado:** Se utilizan para variables categóricas para determinar la asociación entre la variable independiente y la variable dependiente.

30.6.2 Métodos de Wrapper

Los métodos de wrapper evalúan múltiples combinaciones de variables y seleccionan la combinación que optimiza el rendimiento del modelo. Ejemplos incluyen:

- **Selección hacia Adelante:** Comienza con un modelo vacío y agrega variables una por una, seleccionando la variable que mejora más el modelo en cada paso.
- **Selección hacia Atrás:** Comienza con todas las variables y elimina una por una, removiendo la variable que tiene el menor impacto en el modelo en cada paso.
- **Selección Paso a Paso:** Combina la selección hacia adelante y hacia atrás, agregando y eliminando variables según sea necesario.

30.6.3 Métodos Basados en Modelos

Los métodos basados en modelos utilizan técnicas de regularización como Lasso y Ridge para seleccionar variables. Estas técnicas añaden un término de penalización a la función de costo para evitar el sobreajuste.

Regresión Lasso

La regresión Lasso (Least Absolute Shrinkage and Selection Operator) añade una penalización L_1 a la función de costo:

$$J(\beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

donde λ es el parámetro de regularización que controla la cantidad de penalización.

Regresión Ridge

La regresión Ridge añade una penalización L_2 a la función de costo:

$$J(\beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

donde λ es el parámetro de regularización.

30.7 Implementación en R

30.7.1 Limpieza de Datos

Para ilustrar la limpieza de datos en R, considere el siguiente conjunto de datos:

```
data <- data.frame(
  var1 = c(1, 2, 3, NA, 5),
  var2 = c("A", "B", "A", "B", "A"),
  var3 = c(10, 15, 10, 20, 25)
)

# Eliminaci\on de filas con datos faltantes
data_clean <- na.omit(data)

# Imputaci\on de la media
data$var1[is.na(data$var1)] <- mean(data$var1, na.rm = TRUE)
```

30.7.2 Codificación de Variables Categóricas

Para codificar variables categóricas, utilice la función ‘model.matrix’:

```
data <- data.frame(
  var1 = c(1, 2, 3, 4, 5),
  var2 = c("A", "B", "A", "B", "A")
)

# Codificaci\on one-hot
data_onehot <- model.matrix(~ var2 - 1, data = data)
```

30.7.3 Selección de Variables

Para la selección de variables, utilice el paquete ‘caret’:

```
library(caret)

# Dividir los datos en conjuntos de entrenamiento y prueba
set.seed(123)
trainIndex <- createDataPartition(data$var1, p = .8,
                                   list = FALSE,
                                   times = 1)

dataTrain <- data[trainIndex,]
dataTest <- data[-trainIndex,]

# Modelo de regresi\on log\ística
model <- train(var1 ~ ., data = dataTrain, method = "glm", family = "binomial")

# Selecci\on de variables
model <- step(model, direction = "both")
summary(model)
```

Chapter 31

Evaluación del Modelo y Validación Cruzada

31.1 Introducción

Evaluar la calidad y el rendimiento de un modelo de regresión logística es crucial para asegurar que las predicciones sean precisas y útiles. Este capítulo se centra en las técnicas y métricas utilizadas para evaluar modelos de clasificación binaria, así como en la validación cruzada, una técnica para evaluar la generalización del modelo.

31.2 Métricas de Evaluación del Modelo

Las métricas de evaluación permiten cuantificar la precisión y el rendimiento de un modelo. Algunas de las métricas más comunes incluyen:

31.2.1 Curva ROC y AUC

La curva ROC (Receiver Operating Characteristic) es una representación gráfica de la sensibilidad (verdaderos positivos) frente a 1 - especificidad (falsos positivos). El área bajo la curva (AUC) mide la capacidad del modelo para distinguir entre las clases.

$$\begin{aligned}\text{Sensibilidad} &= \frac{\text{TP}}{\text{TP} + \text{FN}} \\ \text{Especificidad} &= \frac{\text{TN}}{\text{TN} + \text{FP}}\end{aligned}$$

31.2.2 Matriz de Confusión

La matriz de confusión es una tabla que muestra el rendimiento del modelo comparando las predicciones con los valores reales. Los términos incluyen:

- **Verdaderos Positivos (TP):** Predicciones correctas de la clase positiva.
- **Falsos Positivos (FP):** Predicciones incorrectas de la clase positiva.

- **Verdaderos Negativos (TN)**: Predicciones correctas de la clase negativa.
- **Falsos Negativos (FN)**: Predicciones incorrectas de la clase negativa.

	Predicción Positiva	Predicción Negativa
Real Positiva	TP	FN
Real Negativa	FP	TN

Table 31.1: Matriz de Confusión

31.2.3 Precisión, Recall y F1-Score

$$\begin{aligned}
 \text{Precisión} &= \frac{TP}{TP + FP} \\
 \text{Recall} &= \frac{TP}{TP + FN} \\
 \text{F1-Score} &= 2 \cdot \frac{\text{Precisión} \cdot \text{Recall}}{\text{Precisión} + \text{Recall}}
 \end{aligned}$$

31.2.4 Log-Loss

La pérdida logarítmica (Log-Loss) mide la precisión de las probabilidades predichas. La fórmula es:

$$\text{Log-Loss} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

donde y_i son los valores reales y p_i son las probabilidades predichas.

31.3 Validación Cruzada

La validación cruzada es una técnica para evaluar la capacidad de generalización de un modelo. Existen varios tipos de validación cruzada:

31.3.1 K-Fold Cross-Validation

En K-Fold Cross-Validation, los datos se dividen en K subconjuntos. El modelo se entrena K veces, cada vez utilizando K-1 subconjuntos para el entrenamiento y el subconjunto restante para la validación.

$$\text{Error Medio} = \frac{1}{K} \sum_{k=1}^K \text{Error}_k$$

31.3.2 Leave-One-Out Cross-Validation (LOOCV)

En LOOCV, cada observación se usa una vez como conjunto de validación y las restantes como conjunto de entrenamiento. Este método es computacionalmente costoso pero útil para conjuntos de datos pequeños.

31.4 Ajuste y Sobreajuste del Modelo

El ajuste adecuado del modelo es crucial para evitar el sobreajuste (overfitting) y el subajuste (underfitting).

31.4.1 Sobreajuste

El sobreajuste ocurre cuando un modelo se ajusta demasiado bien a los datos de entrenamiento, capturando ruido y patrones irrelevantes. Los síntomas incluyen una alta precisión en el entrenamiento y baja precisión en la validación.

31.4.2 Subajuste

El subajuste ocurre cuando un modelo no captura los patrones subyacentes de los datos. Los síntomas incluyen baja precisión tanto en el entrenamiento como en la validación.

31.4.3 Regularización

La regularización es una técnica para prevenir el sobreajuste añadiendo un término de penalización a la función de costo. Las técnicas comunes incluyen:

- Regresión Lasso (L1)
- Regresión Ridge (L2)

31.5 Implementación en R

31.5.1 Evaluación del Modelo

```
# Cargar el paquete necesario
library(caret)

# Dividir los datos en conjuntos de entrenamiento y prueba
set.seed(123)
trainIndex <- createDataPartition(data$var1, p = .8,
                                   list = FALSE,
                                   times = 1)

dataTrain <- data[trainIndex,]
dataTest <- data[-trainIndex,]

# Entrenar el modelo de regresión logística
model <- train(var1 ~ ., data = dataTrain, method = "glm", family = "binomial")

# Predicciones en el conjunto de prueba
predictions <- predict(model, dataTest)

# Matriz de confusión
confusionMatrix(predictions, dataTest$var1)
```

31.5.2 Validación Cruzada

```
# K-Fold Cross-Validation
control <- trainControl(method = "cv", number = 10)
model_cv <- train(var1 ~ ., data = dataTrain, method = "glm",
                  family = "binomial", trControl = control)

# Evaluación del modelo con validación cruzada
print(model_cv)
```

Chapter 32

Diagnóstico del Modelo y Ajuste de Parámetros

32.1 Introducción

El diagnóstico del modelo y el ajuste de parámetros son pasos esenciales para mejorar la precisión y la robustez de los modelos de regresión logística. Este capítulo se enfoca en las técnicas para diagnosticar problemas en los modelos y en métodos para ajustar los parámetros de manera óptima.

32.2 Diagnóstico del Modelo

El diagnóstico del modelo implica evaluar el rendimiento del modelo y detectar posibles problemas, como el sobreajuste, la multicolinealidad y la influencia de puntos de datos individuales.

32.2.1 Residuos

Los residuos son las diferencias entre los valores observados y los valores predichos por el modelo. El análisis de residuos puede revelar patrones que indican problemas con el modelo.

$$\text{Residuo}_i = y_i - \hat{y}_i$$

Residuos Estudiantizados

Los residuos estudiantizados se ajustan por la variabilidad del residuo y se utilizan para detectar outliers.

$$r_i = \frac{\text{Residuo}_i}{\hat{\sigma}\sqrt{1-h_i}}$$

donde h_i es el leverage del punto de datos.

32.2.2 Influencia

La influencia mide el impacto de un punto de datos en los coeficientes del modelo. Los puntos con alta influencia pueden distorsionar el modelo.

Distancia de Cook

La distancia de Cook es una medida de la influencia de un punto de datos en los coeficientes del modelo.

$$D_i = \frac{r_i^2}{p} \cdot \frac{h_i}{1 - h_i}$$

donde p es el número de parámetros en el modelo.

32.2.3 Multicolinealidad

La multicolinealidad ocurre cuando dos o más variables independientes están altamente correlacionadas. Esto puede inflar las varianzas de los coeficientes y hacer que el modelo sea inestable.

Factor de Inflación de la Varianza (VIF)

El VIF mide cuánto se inflan las varianzas de los coeficientes debido a la multicolinealidad.

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

donde R_j^2 es el coeficiente de determinación de la regresión de la variable j contra todas las demás variables.

32.3 Ajuste de Parámetros

El ajuste de parámetros implica seleccionar los valores óptimos para los hiperparámetros del modelo. Esto puede mejorar el rendimiento y prevenir el sobreajuste.

32.3.1 Grid Search

El grid search es un método exhaustivo para ajustar los parámetros. Se define una rejilla de posibles valores de parámetros y se evalúa el rendimiento del modelo para cada combinación.

32.3.2 Random Search

El random search selecciona aleatoriamente combinaciones de valores de parámetros dentro de un rango especificado. Es menos exhaustivo que el grid search, pero puede ser más eficiente.

32.3.3 Bayesian Optimization

La optimización bayesiana utiliza modelos probabilísticos para seleccionar iterativamente los valores de parámetros más prometedores.

32.4 Implementación en R

32.4.1 Diagnóstico del Modelo

```
# Cargar el paquete necesario
library(car)

# Residuos estudentizados
dataTrain$resid <- rstudent(model)
hist(dataTrain$resid, breaks = 20, main = "Residuos Estudentizados")

# Distancia de Cook
dataTrain$cook <- cooks.distance(model)
plot(dataTrain$cook, type = "h", main = "Distancia de Cook")

# Factor de Inflación de la Varianza
vif_values <- vif(model)
print(vif_values)
```

32.4.2 Ajuste de Parámetros

```
# Grid Search con caret
control <- trainControl(method = "cv", number = 10)
tune_grid <- expand.grid(.alpha = c(0, 0.5, 1), .lambda = seq(0.01, 0.1, by = 0.01))

model_tune <- train(var1 ~ ., data = dataTrain, method = "glmnet",
                    trControl = control, tuneGrid = tune_grid)

print(model_tune)
```

Chapter 33

Interpretación de los Resultados

33.1 Introducción

Interpretar correctamente los resultados de un modelo de regresión logística es esencial para tomar decisiones informadas. Este capítulo se centra en la interpretación de los coeficientes del modelo, los odds ratios, los intervalos de confianza y la significancia estadística.

33.2 Coeficientes de Regresión Logística

Los coeficientes de regresión logística representan la relación entre las variables independientes y la variable dependiente en términos de log-odds.

33.2.1 Interpretación de los Coeficientes

Cada coeficiente β_j en el modelo de regresión logística se interpreta como el cambio en el log-odds de la variable dependiente por unidad de cambio en la variable independiente X_j .

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

33.2.2 Signo de los Coeficientes

- **Coeficiente Positivo:** Un coeficiente positivo indica que un aumento en la variable independiente está asociado con un aumento en el log-odds de la variable dependiente.
- **Coeficiente Negativo:** Un coeficiente negativo indica que un aumento en la variable independiente está asociado con una disminución en el log-odds de la variable dependiente.

33.3 Odds Ratios

Las odds ratios proporcionan una interpretación más intuitiva de los coeficientes de regresión logística. La odds ratio para una variable independiente X_j se calcula como e^{β_j} .

33.3.1 Cálculo de las Odds Ratios

$$\text{OR}_j = e^{\beta_j}$$

33.3.2 Interpretación de las Odds Ratios

- **OR > 1:** Un OR mayor que 1 indica que un aumento en la variable independiente está asociado con un aumento en las odds de la variable dependiente.
- **OR < 1:** Un OR menor que 1 indica que un aumento en la variable independiente está asociado con una disminución en las odds de la variable dependiente.
- **OR = 1:** Un OR igual a 1 indica que la variable independiente no tiene efecto sobre las odds de la variable dependiente.

33.4 Intervalos de Confianza

Los intervalos de confianza proporcionan una medida de la incertidumbre asociada con los estimadores de los coeficientes. Un intervalo de confianza del 95% para un coeficiente β_j indica que, en el 95% de las muestras, el intervalo contendrá el valor verdadero de β_j .

33.4.1 Cálculo de los Intervalos de Confianza

Para calcular un intervalo de confianza del 95% para un coeficiente β_j , utilizamos la fórmula:

$$\beta_j \pm 1.96 \cdot \text{SE}(\beta_j)$$

donde $\text{SE}(\beta_j)$ es el error estándar de β_j .

33.5 Significancia Estadística

La significancia estadística se utiliza para determinar si los coeficientes del modelo son significativamente diferentes de cero. Esto se evalúa mediante pruebas de hipótesis.

33.5.1 Prueba de Hipótesis

Para cada coeficiente β_j , la hipótesis nula H_0 es que $\beta_j = 0$. La hipótesis alternativa H_a es que $\beta_j \neq 0$.

33.5.2 P-valor

El p-valor indica la probabilidad de obtener un coeficiente tan extremo como el observado, asumiendo que la hipótesis nula es verdadera. Un p-valor menor que el nivel de significancia α (típicamente 0.05) indica que podemos rechazar la hipótesis nula.

33.6 Implementación en R

33.6.1 Cálculo de Coeficientes y Odds Ratios

```
# Cargar el paquete necesario
library(broom)

# Entrenar el modelo de regresión logística
model <- glm(var1 ~ ., data = dataTrain, family = "binomial")

# Coeficientes del modelo
coef(model)

# Odds ratios
exp(coef(model))
```

33.6.2 Intervalos de Confianza

```
# Intervalos de confianza para los coeficientes
confint(model)

# Intervalos de confianza para las odds ratios
exp(confint(model))
```

33.6.3 P-valores y Significancia Estadística

```
# Resumen del modelo con p-valores
summary(model)
```

Chapter 34

Regresión Logística Multinomial y Análisis de Supervivencia

34.1 Introducción

La regresión logística multinomial y el análisis de supervivencia son extensiones de la regresión logística binaria. Este capítulo se enfoca en las técnicas y aplicaciones de estos métodos avanzados.

34.2 Regresión Logística Multinomial

La regresión logística multinomial se utiliza cuando la variable dependiente tiene más de dos categorías.

34.2.1 Modelo Multinomial

El modelo de regresión logística multinomial generaliza el modelo binario para manejar múltiples categorías. La probabilidad de que una observación pertenezca a la categoría k se expresa como:

$$P(Y = k) = \frac{e^{\beta_{0k} + \beta_{1k}X_1 + \dots + \beta_{nk}X_n}}{\sum_{j=1}^K e^{\beta_{0j} + \beta_{1j}X_1 + \dots + \beta_{nj}X_n}}$$

34.2.2 Estimación de Parámetros

Los coeficientes del modelo multinomial se estiman utilizando máxima verosimilitud, similar a la regresión logística binaria.

34.3 Análisis de Supervivencia

El análisis de supervivencia se utiliza para modelar el tiempo hasta que ocurre un evento de interés, como la muerte o la falla de un componente.

34.3.1 Función de Supervivencia

La función de supervivencia $S(t)$ describe la probabilidad de que una observación sobreviva más allá del tiempo t :

$$S(t) = P(T > t)$$

34.3.2 Modelo de Riesgos Proporcionales de Cox

El modelo de Cox es un modelo de regresión semiparamétrico utilizado para analizar datos de supervivencia:

$$h(t|X) = h_0(t)e^{\beta_1 X_1 + \dots + \beta_p X_p}$$

donde $h(t|X)$ es la tasa de riesgo en el tiempo t dado el vector de covariables X y $h_0(t)$ es la tasa de riesgo basal.

34.4 Implementación en R

34.4.1 Regresión Logística Multinomial

```
# Cargar el paquete necesario
library(nnet)

# Entrenar el modelo de regresión logística multinomial
model_multinom <- multinom(var1 ~ ., data = dataTrain)

# Resumen del modelo
summary(model_multinom)
```

34.4.2 Análisis de Supervivencia

```
# Cargar el paquete necesario
library(survival)

# Crear el objeto de supervivencia
surv_object <- Surv(time = data$time, event = data$status)

# Ajustar el modelo de Cox
model_cox <- coxph(surv_object ~ var1 + var2, data = data)

# Resumen del modelo
summary(model_cox)
```

Chapter 35

Implementación de Regresión Logística en Datos Reales

35.1 Introducción

Implementar un modelo de regresión logística en datos reales implica varias etapas, desde la limpieza de datos hasta la evaluación y validación del modelo. Este capítulo presenta un ejemplo práctico de la implementación de un modelo de regresión logística utilizando un conjunto de datos real.

35.2 Conjunto de Datos

Para este ejemplo, utilizaremos un conjunto de datos disponible públicamente que contiene información sobre clientes bancarios. El objetivo es predecir si un cliente suscribirá un depósito a plazo fijo.

35.3 Preparación de Datos

35.3.1 Carga y Exploración de Datos

Primero, cargamos y exploramos el conjunto de datos para entender su estructura y contenido.

```
# Cargar el paquete necesario
library(dplyr)

# Cargar el conjunto de datos
data <- read.csv("bank.csv")

# Explorar los datos
str(data)
summary(data)
```

35.3.2 Limpieza de Datos

El siguiente paso es limpiar los datos, lo que incluye tratar los valores faltantes y eliminar las duplicidades.


```
# Eliminar duplicados
data <- data %>% distinct()

# Imputar valores faltantes (si existen)
data <- data %>% mutate_if(is.numeric, ~ifelse(is.na(.), mean(., na.rm = TRUE), .))
```

35.3.3 Codificación de Variables Categóricas

Convertimos las variables categóricas en variables numéricas utilizando la codificación one-hot.

```
# Codificación one-hot de variables categóricas
data <- data %>% mutate(across(where(is.factor), ~ as.numeric(as.factor(.))))
```

35.4 División de Datos

Dividimos los datos en conjuntos de entrenamiento y prueba.

```
# Dividir los datos en conjuntos de entrenamiento y prueba
set.seed(123)
trainIndex <- createDataPartition(data$y, p = .8, list = FALSE, times = 1)
dataTrain <- data[trainIndex,]
dataTest <- data[-trainIndex,]
```

35.5 Entrenamiento del Modelo

Entrenamos un modelo de regresión logística utilizando el conjunto de entrenamiento.

```
# Entrenar el modelo de regresión logística
model <- glm(y ~ ., data = dataTrain, family = "binomial")

# Resumen del modelo
summary(model)
```

35.6 Evaluación del Modelo

Evaluamos el rendimiento del modelo utilizando el conjunto de prueba.

```
# Predicciones en el conjunto de prueba
predictions <- predict(model, dataTest, type = "response")

# Convertir probabilidades a etiquetas
predicted_labels <- ifelse(predictions > 0.5, 1, 0)

# Matriz de confusión
confusionMatrix(predicted_labels, dataTest$y)
```

35.7 Interpretación de los Resultados

Interpretamos los coeficientes del modelo y las odds ratios.

```
# Coeficientes del modelo  
coef(model)
```

```
# Odds ratios  
exp(coef(model))
```

Chapter 36

Resumen y Proyecto Final

36.1 Resumen de Conceptos Clave

En este curso, hemos cubierto una variedad de conceptos y técnicas esenciales para la regresión logística. Los conceptos clave incluyen:

- **Fundamentos de Probabilidad y Estadística:** Comprensión de distribuciones de probabilidad, medidas de tendencia central y dispersión, inferencia estadística y pruebas de hipótesis.
- **Regresión Logística:** Modelo de regresión logística binaria y multinomial, interpretación de coeficientes y odds ratios, métodos de estimación y validación.
- **Preparación de Datos:** Limpieza de datos, tratamiento de valores faltantes, codificación de variables categóricas y selección de variables.
- **Evaluación del Modelo:** Curva ROC, AUC, matriz de confusión, precisión, recall, F1-score y validación cruzada.
- **Diagnóstico del Modelo:** Análisis de residuos, influencia, multicolinealidad y ajuste de parámetros.
- **Análisis de Supervivencia:** Modelos de supervivencia, función de supervivencia y modelos de riesgos proporcionales de Cox.

36.2 Buenas Prácticas

Al implementar modelos de regresión logística, es importante seguir buenas prácticas para garantizar la precisión y la robustez de los modelos. Algunas buenas prácticas incluyen:

- **Exploración y Preparación de Datos:** Realizar un análisis exploratorio exhaustivo y preparar los datos adecuadamente antes de construir el modelo.
- **Evaluación y Validación del Modelo:** Utilizar métricas adecuadas para evaluar el rendimiento del modelo y validar el modelo utilizando técnicas como la validación cruzada.
- **Interpretación de Resultados:** Interpretar correctamente los coeficientes del modelo y las odds ratios, y comunicar los resultados de manera clara y concisa.
- **Revisión y Ajuste del Modelo:** Diagnosticar problemas en el modelo y ajustar los parámetros para mejorar el rendimiento.

36.3 Proyecto Final

Para aplicar los conceptos y técnicas aprendidos en este curso, te proponemos realizar un proyecto final utilizando un conjunto de datos de tu elección. El proyecto debe incluir las siguientes etapas:

36.3.1 Selección del Conjunto de Datos

Elige un conjunto de datos relevante que contenga una variable dependiente binaria o multinomial y varias variables independientes.

36.3.2 Exploración y Preparación de Datos

Realiza un análisis exploratorio de los datos y prepara los datos para el modelado. Esto incluye la limpieza de datos, el tratamiento de valores faltantes y la codificación de variables categóricas.

36.3.3 Entrenamiento y Evaluación del Modelo

Entrena un modelo de regresión logística utilizando el conjunto de datos preparado y evalúa su rendimiento utilizando métricas apropiadas.

36.3.4 Interpretación de Resultados

Interpreta los coeficientes del modelo y las odds ratios, y proporciona una explicación clara de los resultados.

36.3.5 Presentación del Proyecto

Presenta tu proyecto en un informe detallado que incluya la descripción del conjunto de datos, los pasos de preparación y modelado, los resultados del modelo y las conclusiones.

Part II

SEGUNDA PARTE: ANALISIS DE SUPERVIVENCIA

Chapter 37

Introducción al Análisis de Supervivencia

37.1 Conceptos Básicos

El análisis de supervivencia es una rama de la estadística que se ocupa del análisis del tiempo que transcurre hasta que ocurre un evento de interés, comúnmente referido como "tiempo de falla". Este campo es ampliamente utilizado en medicina, biología, ingeniería, ciencias sociales, y otros campos.

37.2 Definición de Eventos y Tiempos

En el análisis de supervivencia, un "evento" se refiere a la ocurrencia de un evento específico, como la muerte, la falla de un componente, la recaída de una enfermedad, etc. El "tiempo de supervivencia" es el tiempo que transcurre desde un punto de inicio definido hasta la ocurrencia del evento.

37.3 Censura

La censura ocurre cuando la información completa sobre el tiempo hasta el evento no está disponible para todos los individuos en el estudio. Hay tres tipos principales de censura:

- **Censura a la derecha:** Ocurre cuando el evento de interés no se ha observado para algunos sujetos antes del final del estudio.
- **Censura a la izquierda:** Ocurre cuando el evento de interés ocurrió antes del inicio del periodo de observación.
- **Censura por intervalo:** Ocurre cuando el evento de interés se sabe que ocurrió en un intervalo de tiempo, pero no se conoce el momento exacto.

37.4 Función de Supervivencia

La función de supervivencia, $S(t)$, se define como la probabilidad de que un individuo sobreviva más allá de un tiempo t . Matemáticamente, se expresa como:

$$S(t) = P(T > t)$$

donde T es una variable aleatoria que representa el tiempo hasta el evento. La función de supervivencia tiene las siguientes propiedades:

- $S(0) = 1$: Esto indica que al inicio (tiempo $t = 0$), la probabilidad de haber experimentado el evento es cero, por lo tanto, la supervivencia es del 100
- $\lim_{t \rightarrow \infty} S(t) = 0$: A medida que el tiempo tiende al infinito, la probabilidad de que cualquier individuo aún no haya experimentado el evento tiende a cero.
- $S(t)$ es una función no creciente: Esto significa que a medida que el tiempo avanza, la probabilidad de supervivencia no aumenta.

37.5 Función de Densidad de Probabilidad

La función de densidad de probabilidad $f(t)$ describe la probabilidad de que el evento ocurra en un instante de tiempo específico. Se define como:

$$f(t) = \frac{dF(t)}{dt}$$

donde $F(t)$ es la función de distribución acumulada, $F(t) = P(T \leq t)$. La relación entre $S(t)$ y $f(t)$ es:

$$f(t) = -\frac{dS(t)}{dt}$$

37.6 Función de Riesgo

La función de riesgo, $\lambda(t)$, también conocida como función de tasa de fallas o hazard rate, se define como la tasa instantánea de ocurrencia del evento en el tiempo t , dado que el individuo ha sobrevivido hasta el tiempo t . Matemáticamente, se expresa como:

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t}$$

Esto se puede reescribir usando $f(t)$ y $S(t)$ como:

$$\lambda(t) = \frac{f(t)}{S(t)}$$

37.7 Relación entre Función de Supervivencia y Función de Riesgo

La función de supervivencia y la función de riesgo están relacionadas a través de la siguiente ecuación:

$$S(t) = \exp \left(- \int_0^t \lambda(u) du \right)$$

Esta fórmula se deriva del hecho de que la función de supervivencia es la probabilidad acumulativa de no haber experimentado el evento hasta el tiempo t , y $\lambda(t)$ es la tasa instantánea de ocurrencia del evento.

La función de riesgo también puede ser expresada como:

$$\lambda(t) = -\frac{d}{dt} \log S(t)$$

37.8 Deducción de la Función de Supervivencia

La relación entre la función de supervivencia y la función de riesgo se puede deducir integrando la función de riesgo:

$$\begin{aligned}S(t) &= \exp\left(-\int_0^t \lambda(u) du\right) \\ \log S(t) &= -\int_0^t \lambda(u) du \\ \frac{d}{dt} \log S(t) &= -\lambda(t) \\ \lambda(t) &= -\frac{d}{dt} \log S(t)\end{aligned}$$

37.9 Ejemplo de Cálculo

Supongamos que tenemos una muestra de tiempos de supervivencia T_1, T_2, \dots, T_n . Podemos estimar la función de supervivencia empírica como:

$$\hat{S}(t) = \frac{\text{Número de individuos que sobreviven más allá de } t}{\text{Número total de individuos en riesgo en } t}$$

y la función de riesgo empírica como:

$$\hat{\lambda}(t) = \frac{\text{Número de eventos en } t}{\text{Número de individuos en riesgo en } t}$$

37.10 Conclusión

El análisis de supervivencia es una herramienta poderosa para analizar datos de tiempo hasta evento. Entender los conceptos básicos como la función de supervivencia y la función de riesgo es fundamental para el análisis más avanzado.

Chapter 38

Función de Supervivencia y Función de Riesgo

38.1 Introducción

Este capítulo profundiza en la definición y propiedades de la función de supervivencia y la función de riesgo, dos conceptos fundamentales en el análisis de supervivencia. Entender estas funciones y su relación es crucial para modelar y analizar datos de tiempo hasta evento.

38.2 Función de Supervivencia

La función de supervivencia, $S(t)$, describe la probabilidad de que un individuo sobreviva más allá de un tiempo t . Formalmente, se define como:

$$S(t) = P(T > t)$$

donde T es una variable aleatoria que representa el tiempo hasta el evento.

38.2.1 Propiedades de la Función de Supervivencia

La función de supervivencia tiene varias propiedades importantes:

- $S(0) = 1$: Indica que la probabilidad de haber experimentado el evento en el tiempo 0 es cero.
- $\lim_{t \rightarrow \infty} S(t) = 0$: A medida que el tiempo tiende al infinito, la probabilidad de supervivencia tiende a cero.
- $S(t)$ es una función no creciente: A medida que el tiempo avanza, la probabilidad de supervivencia no aumenta.

38.2.2 Derivación de $S(t)$

Si la función de densidad de probabilidad $f(t)$ del tiempo de supervivencia T es conocida, la función de supervivencia puede derivarse como:

$$\begin{aligned} S(t) &= P(T > t) \\ &= 1 - P(T \leq t) \\ &= 1 - F(t) \\ &= 1 - \int_0^t f(u) du \end{aligned}$$

donde $F(t)$ es la función de distribución acumulada.

38.2.3 Ejemplo de Cálculo de $S(t)$

Consideremos un ejemplo donde el tiempo de supervivencia T sigue una distribución exponencial con tasa λ . La función de densidad de probabilidad $f(t)$ es:

$$f(t) = \lambda e^{-\lambda t}, \quad t \geq 0$$

La función de distribución acumulada $F(t)$ es:

$$F(t) = \int_0^t \lambda e^{-\lambda u} du = 1 - e^{-\lambda t}$$

Por lo tanto, la función de supervivencia $S(t)$ es:

$$S(t) = 1 - F(t) = e^{-\lambda t}$$

38.3 Función de Riesgo

La función de riesgo, $\lambda(t)$, proporciona la tasa instantánea de ocurrencia del evento en el tiempo t , dado que el individuo ha sobrevivido hasta el tiempo t . Matemáticamente, se define como:

$$\lambda(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t}$$

38.3.1 Relación entre $\lambda(t)$ y $f(t)$

La función de riesgo se puede relacionar con la función de densidad de probabilidad $f(t)$ y la función de supervivencia $S(t)$ de la siguiente manera:

$$\lambda(t) = \frac{f(t)}{S(t)}$$

38.3.2 Derivación de $\lambda(t)$

La derivación de $\lambda(t)$ se basa en la definición condicional de la probabilidad:

$$\begin{aligned}\lambda(t) &= \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t} \\&= \lim_{\Delta t \rightarrow 0} \frac{\frac{P(t \leq T < t + \Delta t \text{ y } T \geq t)}{P(T \geq t)}}{\Delta t} \\&= \lim_{\Delta t \rightarrow 0} \frac{\frac{P(t \leq T < t + \Delta t)}{P(T \geq t)}}{\Delta t} \\&= \frac{f(t)}{S(t)}\end{aligned}$$

38.4 Relación entre Función de Supervivencia y Función de Riesgo

La función de supervivencia y la función de riesgo están estrechamente relacionadas. La relación se expresa mediante la siguiente ecuación:

$$S(t) = \exp \left(- \int_0^t \lambda(u) du \right)$$

38.4.1 Deducción de la Relación

Para deducir esta relación, consideramos la derivada logarítmica de la función de supervivencia:

$$\begin{aligned}S(t) &= \exp \left(- \int_0^t \lambda(u) du \right) \\ \log S(t) &= - \int_0^t \lambda(u) du \\ \frac{d}{dt} \log S(t) &= -\lambda(t) \\ \lambda(t) &= -\frac{d}{dt} \log S(t)\end{aligned}$$

38.5 Interpretación de la Función de Riesgo

La función de riesgo, $\lambda(t)$, se interpreta como la tasa instantánea de ocurrencia del evento por unidad de tiempo, dado que el individuo ha sobrevivido hasta el tiempo t . Es una medida local del riesgo de falla en un instante específico.

38.5.1 Ejemplo de Cálculo de $\lambda(t)$

Consideremos nuevamente el caso donde el tiempo de supervivencia T sigue una distribución exponencial con tasa λ . La función de densidad de probabilidad $f(t)$ es:

$$f(t) = \lambda e^{-\lambda t}$$

La función de supervivencia $S(t)$ es:

$$S(t) = e^{-\lambda t}$$

La función de riesgo $\lambda(t)$ se calcula como:

$$\begin{aligned}\lambda(t) &= \frac{f(t)}{S(t)} \\ &= \frac{\lambda e^{-\lambda t}}{e^{-\lambda t}} \\ &= \lambda\end{aligned}$$

En este caso, $\lambda(t)$ es constante y igual a λ , lo que es una característica de la distribución exponencial.

38.6 Funciones de Riesgo Acumulada y Media Residual

La función de riesgo acumulada $H(t)$ se define como:

$$H(t) = \int_0^t \lambda(u) du$$

Esta función proporciona la suma acumulada de la tasa de riesgo hasta el tiempo t .

La función de vida media residual $e(t)$ se define como la esperanza del tiempo de vida restante dado que el individuo ha sobrevivido hasta el tiempo t :

$$e(t) = \mathbb{E}[T - t \mid T > t] = \int_t^\infty S(u) du$$

38.7 Ejemplo de Cálculo de Función de Riesgo Acumulada y Vida Media Residual

Consideremos nuevamente la distribución exponencial con tasa λ . La función de riesgo acumulada $H(t)$ es:

$$\begin{aligned}H(t) &= \int_0^t \lambda du \\ &= \lambda t\end{aligned}$$

La función de vida media residual $e(t)$ es:

$$\begin{aligned}e(t) &= \int_t^\infty e^{-\lambda u} du \\ &= \left[\frac{-1}{\lambda} e^{-\lambda u} \right]_t^\infty \\ &= \frac{1}{\lambda} e^{-\lambda t} \\ &= \frac{1}{\lambda}\end{aligned}$$

En este caso, la vida media residual es constante e igual a $\frac{1}{\lambda}$, otra característica de la distribución exponencial.

38.8 Conclusión

La función de supervivencia y la función de riesgo son herramientas fundamentales en el análisis de supervivencia. Entender su definición, propiedades, y la relación entre ellas es esencial para modelar y analizar correctamente los datos de tiempo hasta evento. Las funciones de riesgo acumulada y vida media residual proporcionan información adicional sobre la dinámica del riesgo a lo largo del tiempo.

Chapter 39

Estimador de Kaplan-Meier

39.1 Introducción

El estimador de Kaplan-Meier, también conocido como la función de supervivencia empírica, es una herramienta no paramétrica para estimar la función de supervivencia a partir de datos censurados. Este método es especialmente útil cuando los tiempos de evento están censurados a la derecha.

39.2 Definición del Estimador de Kaplan-Meier

El estimador de Kaplan-Meier se define como:

$$\hat{S}(t) = \prod_{t_i \leq t} \left(1 - \frac{d_i}{n_i}\right)$$

donde:

- t_i es el tiempo del i -ésimo evento,
- d_i es el número de eventos que ocurren en t_i ,
- n_i es el número de individuos en riesgo justo antes de t_i .

39.3 Propiedades del Estimador de Kaplan-Meier

El estimador de Kaplan-Meier tiene las siguientes propiedades:

- Es una función escalonada que disminuye en los tiempos de los eventos observados.
- Puede manejar datos censurados a la derecha.
- Proporciona una estimación no paramétrica de la función de supervivencia.

39.3.1 Función Escalonada

La función escalonada del estimador de Kaplan-Meier significa que $\hat{S}(t)$ permanece constante entre los tiempos de los eventos y disminuye en los tiempos de los eventos. Matemáticamente, si t_i es el tiempo del i -ésimo evento, entonces:

$$\hat{S}(t) = \hat{S}(t_i) \quad \text{para } t_i \leq t < t_{i+1}$$

39.3.2 Manejo de Datos Censurados

El estimador de Kaplan-Meier maneja datos censurados a la derecha al ajustar la estimación de la función de supervivencia sólo en los tiempos en que ocurren eventos. Si un individuo es censurado antes de experimentar el evento, no contribuye a la disminución de $\hat{S}(t)$ en el tiempo de censura. Esto asegura que la censura no sesga la estimación de la supervivencia.

39.3.3 Estimación No Paramétrica

El estimador de Kaplan-Meier es no paramétrico porque no asume ninguna forma específica para la distribución de los tiempos de supervivencia. En cambio, utiliza la información empírica disponible para estimar la función de supervivencia.

39.4 Deducción del Estimador de Kaplan-Meier

La deducción del estimador de Kaplan-Meier se basa en el principio de probabilidad condicional. Consideremos un conjunto de tiempos de supervivencia observados t_1, t_2, \dots, t_k con eventos en cada uno de estos tiempos. El estimador de la probabilidad de supervivencia más allá del tiempo t es el producto de las probabilidades de sobrevivir más allá de cada uno de los tiempos de evento observados hasta t .

39.4.1 Probabilidad Condicional

La probabilidad de sobrevivir más allá de t_i , dado que el individuo ha sobrevivido justo antes de t_i , es:

$$P(T > t_i \mid T \geq t_i) = 1 - \frac{d_i}{n_i}$$

donde d_i es el número de eventos en t_i y n_i es el número de individuos en riesgo justo antes de t_i .

39.4.2 Producto de Probabilidades Condicionales

La probabilidad de sobrevivir más allá de un tiempo t cualquiera, dada la secuencia de tiempos de evento, es el producto de las probabilidades condicionales de sobrevivir más allá de cada uno de los tiempos de evento observados hasta t . Así, el estimador de Kaplan-Meier se obtiene como:

$$\hat{S}(t) = \prod_{t_i \leq t} \left(1 - \frac{d_i}{n_i}\right)$$

39.5 Ejemplo de Cálculo

Supongamos que tenemos los siguientes tiempos de supervivencia observados para cinco individuos: 2, 3, 5, 7, 8. Supongamos además que tenemos censura a la derecha en el tiempo 10. Los tiempos de evento y el número de individuos en riesgo justo antes de cada evento son:

Tiempo (t_i)	Eventos (d_i)	En Riesgo (n_i)
2	1	5
3	1	4
5	1	3
7	1	2
8	1	1

Table 39.1: Ejemplo de cálculo del estimador de Kaplan-Meier

Usando estos datos, el estimador de Kaplan-Meier se calcula como:

$$\begin{aligned}
\hat{S}(2) &= 1 - \frac{1}{5} = 0.8 \\
\hat{S}(3) &= 0.8 \times \left(1 - \frac{1}{4}\right) = 0.8 \times 0.75 = 0.6 \\
\hat{S}(5) &= 0.6 \times \left(1 - \frac{1}{3}\right) = 0.6 \times 0.6667 = 0.4 \\
\hat{S}(7) &= 0.4 \times \left(1 - \frac{1}{2}\right) = 0.4 \times 0.5 = 0.2 \\
\hat{S}(8) &= 0.2 \times \left(1 - \frac{1}{1}\right) = 0.2 \times 0 = 0
\end{aligned}$$

39.6 Intervalos de Confianza para el Estimador de Kaplan-Meier

Para calcular intervalos de confianza para el estimador de Kaplan-Meier, se puede usar la transformación logarítmica y la aproximación normal. Un intervalo de confianza aproximado para $\log(-\log(\hat{S}(t)))$ se obtiene como:

$$\log(-\log(\hat{S}(t))) \pm z_{\alpha/2} \sqrt{\frac{1}{d_i(n_i - d_i)}}$$

donde $z_{\alpha/2}$ es el percentil correspondiente de la distribución normal estándar.

39.7 Transformación Logarítmica Inversa

La transformación logarítmica inversa se utiliza para obtener los límites del intervalo de confianza para $S(t)$:

$$\hat{S}(t) = \exp \left(-\exp \left(\log(-\log(\hat{S}(t))) \pm z_{\alpha/2} \sqrt{\frac{1}{d_i(n_i - d_i)}} \right) \right)$$

39.8 Cálculo Detallado de Intervalos de Confianza

Para un cálculo más detallado de los intervalos de confianza, consideremos un tiempo específico t_j . La varianza del estimador de Kaplan-Meier en t_j se puede estimar usando Greenwood's formula:

$$\text{Var}(\hat{S}(t_j)) = \hat{S}(t_j)^2 \sum_{t_i \leq t_j} \frac{d_i}{n_i(n_i - d_i)}$$

El intervalo de confianza aproximado para $\hat{S}(t_j)$ es entonces:

$$\hat{S}(t_j) \pm z_{\alpha/2} \sqrt{\text{Var}(\hat{S}(t_j))}$$

39.9 Ejemplo de Intervalo de Confianza

Supongamos que en el ejemplo anterior queremos calcular el intervalo de confianza para $\hat{S}(3)$. Primero, calculamos la varianza:

$$\begin{aligned} \text{Var}(\hat{S}(3)) &= \hat{S}(3)^2 \left(\frac{1}{5 \times 4} + \frac{1}{4 \times 3} \right) \\ &= 0.6^2 \left(\frac{1}{20} + \frac{1}{12} \right) \\ &= 0.36 (0.05 + 0.0833) \\ &= 0.36 \times 0.1333 \\ &= 0.048 \end{aligned}$$

El intervalo de confianza es entonces:

$$0.6 \pm 1.96\sqrt{0.048} = 0.6 \pm 1.96 \times 0.219 = 0.6 \pm 0.429$$

Por lo tanto, el intervalo de confianza para $\hat{S}(3)$ es aproximadamente (0.171, 1.029). Dado que una probabilidad no puede exceder 1, ajustamos el intervalo a (0.171, 1.0).

39.10 Interpretación del Estimador de Kaplan-Meier

El estimador de Kaplan-Meier proporciona una estimación empírica de la función de supervivencia que es fácil de interpretar y calcular. Su capacidad para manejar datos censurados lo hace especialmente útil en estudios de supervivencia.

39.11 Conclusión

El estimador de Kaplan-Meier es una herramienta poderosa para estimar la función de supervivencia en presencia de datos censurados. Su cálculo es relativamente sencillo y proporciona una estimación no paramétrica robusta de la supervivencia a lo largo del tiempo. La interpretación adecuada de este estimador y su intervalo de confianza asociado es fundamental para el análisis de datos de supervivencia.

Chapter 40

Comparación de Curvas de Supervivencia

40.1 Introducción

Comparar curvas de supervivencia es crucial para determinar si existen diferencias significativas en las tasas de supervivencia entre diferentes grupos. Las pruebas de hipótesis, como el test de log-rank, son herramientas comunes para esta comparación.

40.2 Test de Log-rank

El test de log-rank se utiliza para comparar las curvas de supervivencia de dos o más grupos. La hipótesis nula es que no hay diferencia en las funciones de riesgo entre los grupos.

40.2.1 Fórmula del Test de Log-rank

El estadístico del test de log-rank se define como:

$$\chi^2 = \frac{\left(\sum_{i=1}^k (O_i - E_i)\right)^2}{\sum_{i=1}^k V_i}$$

donde:

- O_i es el número observado de eventos en el grupo i .
- E_i es el número esperado de eventos en el grupo i .
- V_i es la varianza del número de eventos en el grupo i .

40.2.2 Cálculo de E_i y V_i

El número esperado de eventos E_i y la varianza V_i se calculan como:

$$\begin{aligned} E_i &= \frac{d_i \cdot n_i}{n} \\ V_i &= \frac{d_i \cdot (n - d_i) \cdot n_i \cdot (n - n_i)}{n^2 \cdot (n - 1)} \end{aligned}$$

donde:

- d_i es el número total de eventos en el grupo i .
- n_i es el número de individuos en riesgo en el grupo i .
- n es el número total de individuos en todos los grupos.

40.3 Ejemplo de Cálculo del Test de Log-rank

Supongamos que tenemos dos grupos con los siguientes datos de eventos:

Grupo	Tiempo (t_i)	Eventos (O_i)	En Riesgo (n_i)
1	2	1	5
1	4	1	4
2	3	1	4
2	5	1	3

Table 40.1: Ejemplo de datos para el test de log-rank

Calculemos E_i y V_i para cada grupo:

$$\begin{aligned}
 E_1 &= \frac{2 \cdot 5}{9} + \frac{2 \cdot 4}{8} = \frac{10}{9} + \frac{8}{8} = 1.11 + 1 = 2.11 \\
 V_1 &= \frac{2 \cdot 7 \cdot 5 \cdot 4}{81 \cdot 8} = \frac{2 \cdot 7 \cdot 5 \cdot 4}{648} = \frac{280}{648} = 0.432 \\
 E_2 &= \frac{2 \cdot 4}{9} + \frac{2 \cdot 3}{8} = \frac{8}{9} + \frac{6}{8} = 0.89 + 0.75 = 1.64 \\
 V_2 &= \frac{2 \cdot 7 \cdot 4 \cdot 4}{81 \cdot 8} = \frac{2 \cdot 7 \cdot 4 \cdot 4}{648} = \frac{224}{648} = 0.346
 \end{aligned}$$

El estadístico de log-rank se calcula como:

$$\begin{aligned}
 \chi^2 &= \frac{((1 - 2.11) + (1 - 1.64))^2}{0.432 + 0.346} \\
 &= \frac{(-1.11 - 0.64)^2}{0.778} \\
 &= \frac{3.04}{0.778} \\
 &= 3.91
 \end{aligned}$$

El valor p se puede obtener comparando χ^2 con una distribución χ^2 con un grado de libertad (dado que estamos comparando dos grupos).

40.4 Interpretación del Test de Log-rank

Un valor p pequeño (generalmente menos de 0.05) indica que hay una diferencia significativa en las curvas de supervivencia entre los grupos. Un valor p grande sugiere que no hay suficiente evidencia para rechazar la hipótesis nula de que las curvas de supervivencia son iguales.

40.5 Pruebas Alternativas

Además del test de log-rank, existen otras pruebas para comparar curvas de supervivencia, como el test de Wilcoxon (Breslow), que da más peso a los eventos en tiempos tempranos.

40.6 Conclusión

El test de log-rank es una herramienta esencial para comparar curvas de supervivencia entre diferentes grupos. Su cálculo se basa en la diferencia entre los eventos observados y esperados en cada grupo, y su interpretación puede ayudar a identificar diferencias significativas en la supervivencia.

Chapter 41

Modelos de Riesgos Proporcionales de Cox

41.1 Introducción

El modelo de riesgos proporcionales de Cox, propuesto por David Cox en 1972, es una de las herramientas más utilizadas en el análisis de supervivencia. Este modelo permite evaluar el efecto de varias covariables en el tiempo hasta el evento, sin asumir una forma específica para la distribución de los tiempos de supervivencia.

41.2 Definición del Modelo de Cox

El modelo de Cox se define como:

$$\lambda(t | X) = \lambda_0(t) \exp(\beta^T X)$$

donde:

- $\lambda(t | X)$ es la función de riesgo en el tiempo t dado el vector de covariables X .
- $\lambda_0(t)$ es la función de riesgo basal en el tiempo t .
- β es el vector de coeficientes del modelo.
- X es el vector de covariables.

41.3 Supuesto de Proporcionalidad de Riesgos

El modelo de Cox asume que las razones de riesgo entre dos individuos son constantes a lo largo del tiempo. Matemáticamente, si X_i y X_j son las covariables de dos individuos, la razón de riesgos se expresa como:

$$\frac{\lambda(t | X_i)}{\lambda(t | X_j)} = \frac{\lambda_0(t) \exp(\beta^T X_i)}{\lambda_0(t) \exp(\beta^T X_j)} = \exp(\beta^T (X_i - X_j))$$

41.4 Estimación de los Parámetros

Los parámetros β se estiman utilizando el método de máxima verosimilitud parcial. La función de verosimilitud parcial se define como:

$$L(\beta) = \prod_{i=1}^k \frac{\exp(\beta^T X_i)}{\sum_{j \in R(t_i)} \exp(\beta^T X_j)}$$

donde $R(t_i)$ es el conjunto de individuos en riesgo en el tiempo t_i .

41.4.1 Función de Log-Verosimilitud Parcial

La función de log-verosimilitud parcial es:

$$\log L(\beta) = \sum_{i=1}^k \left(\beta^T X_i - \log \sum_{j \in R(t_i)} \exp(\beta^T X_j) \right)$$

41.4.2 Derivadas Parciales y Maximización

Para encontrar los estimadores de máxima verosimilitud, resolvemos el sistema de ecuaciones obtenido al igualar a cero las derivadas parciales de $\log L(\beta)$ con respecto a β :

$$\frac{\partial \log L(\beta)}{\partial \beta} = \sum_{i=1}^k \left(X_i - \frac{\sum_{j \in R(t_i)} X_j \exp(\beta^T X_j)}{\sum_{j \in R(t_i)} \exp(\beta^T X_j)} \right) = 0$$

41.5 Interpretación de los Coeficientes

Cada coeficiente β_i representa el logaritmo de la razón de riesgos asociado con un incremento unitario en la covariable X_i . Un valor positivo de β_i indica que un aumento en X_i incrementa el riesgo del evento, mientras que un valor negativo indica una reducción del riesgo.

41.6 Evaluación del Modelo

El modelo de Cox se evalúa utilizando varias técnicas, como el análisis de residuos de Schoenfeld para verificar el supuesto de proporcionalidad de riesgos, y el uso de curvas de supervivencia estimadas para evaluar la bondad de ajuste.

41.6.1 Residuos de Schoenfeld

Los residuos de Schoenfeld se utilizan para evaluar la proporcionalidad de riesgos. Para cada evento en el tiempo t_i , el residuo de Schoenfeld para la covariable X_j se define como:

$$r_{ij} = X_{ij} - \hat{X}_{ij}$$

donde \hat{X}_{ij} es la covariable ajustada.

41.6.2 Curvas de Supervivencia Ajustadas

Las curvas de supervivencia ajustadas se obtienen utilizando la función de riesgo basal estimada y los coeficientes del modelo. La función de supervivencia ajustada se define como:

$$\hat{S}(t | X) = \hat{S}_0(t)^{\exp(\beta^T X)}$$

donde $\hat{S}_0(t)$ es la función de supervivencia basal estimada.

41.7 Ejemplo de Aplicación del Modelo de Cox

Consideremos un ejemplo con tres covariables: edad, sexo y tratamiento. Supongamos que los datos se ajustan a un modelo de Cox y obtenemos los siguientes coeficientes:

$$\hat{\beta}_{edad} = 0.02, \quad \hat{\beta}_{sexo} = -0.5, \quad \hat{\beta}_{tratamiento} = 1.2$$

La función de riesgo ajustada se expresa como:

$$\lambda(t | X) = \lambda_0(t) \exp(0.02 \cdot edad - 0.5 \cdot sexo + 1.2 \cdot tratamiento)$$

41.8 Conclusión

El modelo de riesgos proporcionales de Cox es una herramienta poderosa para analizar datos de supervivencia con múltiples covariables. Su flexibilidad y la falta de suposiciones fuertes sobre la distribución de los tiempos de supervivencia lo hacen ampliamente aplicable en diversas disciplinas.

Chapter 42

Diagnóstico y Validación de Modelos de Cox

42.1 Introducción

Una vez ajustado un modelo de Cox, es crucial realizar diagnósticos y validaciones para asegurar que el modelo es apropiado y que los supuestos subyacentes son válidos. Esto incluye la verificación del supuesto de proporcionalidad de riesgos y la evaluación del ajuste del modelo.

42.2 Supuesto de Proporcionalidad de Riesgos

El supuesto de proporcionalidad de riesgos implica que la razón de riesgos entre dos individuos es constante a lo largo del tiempo. Si este supuesto no se cumple, las inferencias hechas a partir del modelo pueden ser incorrectas.

42.2.1 Residuos de Schoenfeld

Los residuos de Schoenfeld se utilizan para evaluar la proporcionalidad de riesgos. Para cada evento en el tiempo t_i , el residuo de Schoenfeld para la covariable X_j se define como:

$$r_{ij} = X_{ij} - \hat{X}_{ij}$$

donde \hat{X}_{ij} es la covariable ajustada. Si los residuos de Schoenfeld no muestran una tendencia sistemática cuando se trazan contra el tiempo, el supuesto de proporcionalidad de riesgos es razonable.

42.3 Bondad de Ajuste

La bondad de ajuste del modelo de Cox se evalúa comparando las curvas de supervivencia observadas y ajustadas, y utilizando estadísticas de ajuste global.

42.3.1 Curvas de Supervivencia Ajustadas

Las curvas de supervivencia ajustadas se obtienen utilizando la función de riesgo basal estimada y los coeficientes del modelo. La función de supervivencia ajustada se define como:

$$\hat{S}(t | X) = \hat{S}_0(t)^{\exp(\beta^T X)}$$

donde $\hat{S}_0(t)$ es la función de supervivencia basal estimada. Comparar estas curvas con las curvas de Kaplan-Meier para diferentes niveles de las covariables puede proporcionar una validación visual del ajuste del modelo.

42.3.2 Estadísticas de Ajuste Global

Las estadísticas de ajuste global, como el test de la desviación y el test de la bondad de ajuste de Grambsch y Therneau, se utilizan para evaluar el ajuste global del modelo de Cox.

42.4 Diagnóstico de Influencia

El diagnóstico de influencia identifica observaciones individuales que tienen un gran impacto en los estimados del modelo. Los residuos de devianza y los residuos de martingala se utilizan comúnmente para este propósito.

42.4.1 Residuos de Deviance

Los residuos de deviance se definen como:

$$D_i = \text{sign}(O_i - E_i) \sqrt{-2 \left(O_i \log \frac{O_i}{E_i} - (O_i - E_i) \right)}$$

donde O_i es el número observado de eventos y E_i es el número esperado de eventos. Observaciones con residuos de deviance grandes en valor absoluto pueden ser influyentes.

42.4.2 Residuos de Martingala

Los residuos de martingala se definen como:

$$M_i = O_i - E_i$$

donde O_i es el número observado de eventos y E_i es el número esperado de eventos. Los residuos de martingala se utilizan para detectar observaciones que no se ajustan bien al modelo.

42.5 Ejemplo de Diagnóstico

Consideremos un modelo de Cox ajustado con las covariables edad, sexo y tratamiento. Para diagnosticar la influencia de observaciones individuales, calculamos los residuos de deviance y martingala para cada observación.

Observaciones con residuos de deviance grandes en valor absoluto (como la observación 3) pueden ser influyentes y requieren una revisión adicional.

Observación	Edad	Sexo	Tratamiento	Residuo de Deviance
1	50	0	1	1.2
2	60	1	0	-0.5
3	45	0	1	-1.8
4	70	1	0	0.3

Table 42.1: Residuos de deviance para observaciones individuales

42.6 Conclusión

El diagnóstico y la validación son pasos críticos en el análisis de modelos de Cox. Evaluar el supuesto de proporcionalidad de riesgos, la bondad de ajuste y la influencia de observaciones individuales asegura que las inferencias y conclusiones derivadas del modelo sean válidas y fiables.

Chapter 43

Modelos Acelerados de Fallos

43.1 Introducción

Los modelos acelerados de fallos (AFT) son una alternativa a los modelos de riesgos proporcionales de Cox. En lugar de asumir que las covariables afectan la tasa de riesgo, los modelos AFT asumen que las covariables multiplican el tiempo de supervivencia por una constante.

43.2 Definición del Modelo AFT

Un modelo AFT se expresa como:

$$T = T_0 \exp(\beta^T X)$$

donde:

- T es el tiempo de supervivencia observado.
- T_0 es el tiempo de supervivencia bajo condiciones basales.
- β es el vector de coeficientes del modelo.
- X es el vector de covariables.

43.2.1 Transformación Logarítmica

El modelo AFT se puede transformar logarítmicamente para obtener una forma lineal:

$$\log(T) = \log(T_0) + \beta^T X$$

43.3 Estimación de los Parámetros

Los parámetros del modelo AFT se estiman utilizando el método de máxima verosimilitud. La función de verosimilitud se define como:

$$L(\beta) = \prod_{i=1}^n f(t_i | X_i; \beta)$$

donde $f(t_i | X_i; \beta)$ es la función de densidad de probabilidad del tiempo de supervivencia t_i dado el vector de covariables X_i y los parámetros β .

43.3.1 Función de Log-Verosimilitud

La función de log-verosimilitud es:

$$\log L(\beta) = \sum_{i=1}^n \log f(t_i | X_i; \beta)$$

43.3.2 Maximización de la Verosimilitud

Los estimadores de máxima verosimilitud se obtienen resolviendo el sistema de ecuaciones obtenido al igualar a cero las derivadas parciales de $\log L(\beta)$ con respecto a β :

$$\frac{\partial \log L(\beta)}{\partial \beta} = 0$$

43.4 Distribuciones Comunes en Modelos AFT

En los modelos AFT, el tiempo de supervivencia T puede seguir varias distribuciones comunes, como la exponencial, Weibull, log-normal y log-logística. Cada una de estas distribuciones tiene diferentes propiedades y aplicaciones.

43.4.1 Modelo Exponencial AFT

En un modelo exponencial AFT, el tiempo de supervivencia T sigue una distribución exponencial con parámetro λ :

$$f(t) = \lambda \exp(-\lambda t)$$

La función de supervivencia es:

$$S(t) = \exp(-\lambda t)$$

La transformación logarítmica del tiempo de supervivencia es:

$$\log(T) = \log\left(\frac{1}{\lambda}\right) + \beta^T X$$

43.4.2 Modelo Weibull AFT

En un modelo Weibull AFT, el tiempo de supervivencia T sigue una distribución Weibull con parámetros λ y k :

$$f(t) = \lambda k t^{k-1} \exp(-\lambda t^k)$$

La función de supervivencia es:

$$S(t) = \exp(-\lambda t^k)$$

La transformación logarítmica del tiempo de supervivencia es:

$$\log(T) = \log\left(\left(\frac{1}{\lambda}\right)^{1/k}\right) + \frac{\beta^T X}{k}$$

43.5 Interpretación de los Coeficientes

En los modelos AFT, los coeficientes β_i se interpretan como factores multiplicativos del tiempo de supervivencia. Un valor positivo de β_i indica que un aumento en la covariable X_i incrementa el tiempo de supervivencia, mientras que un valor negativo indica una reducción del tiempo de supervivencia.

43.6 Ejemplo de Aplicación del Modelo AFT

Consideremos un ejemplo con tres covariables: edad, sexo y tratamiento. Supongamos que los datos se ajustan a un modelo Weibull AFT y obtenemos los siguientes coeficientes:

$$\hat{\beta}_{edad} = -0.02, \quad \hat{\beta}_{sexo} = 0.5, \quad \hat{\beta}_{tratamiento} = -1.2$$

La función de supervivencia ajustada se expresa como:

$$S(t | X) = \exp \left(- \left(\frac{t \exp(-0.02 \cdot edad + 0.5 \cdot sexo - 1.2 \cdot tratamiento)}{\lambda} \right)^k \right)$$

43.7 Conclusión

Los modelos AFT proporcionan una alternativa flexible a los modelos de riesgos proporcionales de Cox. Su enfoque en la multiplicación del tiempo de supervivencia por una constante permite una interpretación intuitiva y aplicaciones en diversas áreas.

Chapter 44

Análisis Multivariado de Supervivencia

44.1 Introducción

El análisis multivariado de supervivencia extiende los modelos de supervivencia para incluir múltiples covariables, permitiendo evaluar su efecto simultáneo sobre el tiempo hasta el evento. Los modelos de Cox y AFT son comúnmente utilizados en este contexto.

44.2 Modelo de Cox Multivariado

El modelo de Cox multivariado se define como:

$$\lambda(t | X) = \lambda_0(t) \exp(\beta^T X)$$

donde X es un vector de covariables.

44.2.1 Estimación de los Parámetros

Los parámetros β se estiman utilizando el método de máxima verosimilitud parcial, como se discutió anteriormente. La función de verosimilitud parcial se maximiza para obtener los estimadores de los coeficientes.

44.3 Modelo AFT Multivariado

El modelo AFT multivariado se expresa como:

$$T = T_0 \exp(\beta^T X)$$

44.3.1 Estimación de los Parámetros

Los parámetros β se estiman utilizando el método de máxima verosimilitud, similar al caso univariado. La función de verosimilitud se maximiza para obtener los estimadores de los coeficientes.

44.4 Interacción y Efectos No Lineales

En el análisis multivariado, es importante considerar la posibilidad de interacciones entre covariables y efectos no lineales. Estos se pueden incluir en los modelos extendiendo las funciones de riesgo o supervivencia.

44.4.1 Interacciones

Las interacciones entre covariables se pueden modelar añadiendo términos de interacción en el modelo:

$$\lambda(t | X) = \lambda_0(t) \exp(\beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2)$$

donde $X_1 X_2$ es el término de interacción.

44.4.2 Efectos No Lineales

Los efectos no lineales se pueden modelar utilizando funciones no lineales de las covariables, como polinomios o splines:

$$\lambda(t | X) = \lambda_0(t) \exp(\beta_1 X + \beta_2 X^2)$$

44.5 Selección de Variables

La selección de variables es crucial en el análisis multivariado para evitar el sobreajuste y mejorar la interpretabilidad del modelo. Métodos como la regresión hacia atrás, la regresión hacia adelante y la selección por criterios de información (AIC, BIC) son comúnmente utilizados.

44.5.1 Regresión Hacia Atrás

La regresión hacia atrás comienza con todas las covariables en el modelo y elimina iterativamente la covariable menos significativa hasta que todas las covariables restantes sean significativas.

44.5.2 Regresión Hacia Adelante

La regresión hacia adelante comienza con un modelo vacío y añade iterativamente la covariable más significativa hasta que no se pueda añadir ninguna covariable adicional significativa.

44.5.3 Criterios de Información

Los criterios de información, como el AIC (Akaike Information Criterion) y el BIC (Bayesian Information Criterion), se utilizan para seleccionar el modelo que mejor se ajusta a los datos con la menor complejidad posible:

$$\begin{aligned} AIC &= -2 \log L + 2k \\ BIC &= -2 \log L + k \log n \end{aligned}$$

donde L es la función de verosimilitud del modelo, k es el número de parámetros en el modelo y n es el tamaño de la muestra.

44.6 Ejemplo de Análisis Multivariado

Consideremos un ejemplo con tres covariables: edad, sexo y tratamiento. Ajustamos un modelo de Cox multivariado y obtenemos los siguientes coeficientes:

$$\hat{\beta}_{edad} = 0.03, \quad \hat{\beta}_{sexo} = -0.6, \quad \hat{\beta}_{tratamiento} = 1.5$$

La función de riesgo ajustada se expresa como:

$$\lambda(t | X) = \lambda_0(t) \exp(0.03 \cdot edad - 0.6 \cdot sexo + 1.5 \cdot tratamiento)$$

44.7 Conclusión

El análisis multivariado de supervivencia permite evaluar el efecto conjunto de múltiples covariables sobre el tiempo hasta el evento. La inclusión de interacciones y efectos no lineales, junto con la selección adecuada de variables, mejora la precisión y la interpretabilidad de los modelos de supervivencia.

Chapter 45

Supervivencia en Datos Complicados

45.1 Introducción

El análisis de supervivencia en datos complicados se refiere a la evaluación de datos de supervivencia que presentan desafíos adicionales, como la censura por intervalo, datos truncados y datos con múltiples tipos de eventos. Estos escenarios requieren métodos avanzados para un análisis adecuado.

45.2 Censura por Intervalo

La censura por intervalo ocurre cuando el evento de interés se sabe que ocurrió dentro de un intervalo de tiempo, pero no se conoce el momento exacto. Esto es común en estudios donde las observaciones se realizan en puntos de tiempo discretos.

45.2.1 Modelo para Datos Censurados por Intervalo

Para datos censurados por intervalo, la función de verosimilitud se modifica para incluir la probabilidad de que el evento ocurra dentro de un intervalo:

$$L(\beta) = \prod_{i=1}^n P(T_i \in [L_i, U_i] \mid X_i; \beta)$$

donde $[L_i, U_i]$ es el intervalo de tiempo durante el cual se sabe que ocurrió el evento para el individuo i .

45.3 Datos Truncados

Los datos truncados ocurren cuando los tiempos de supervivencia están sujetos a un umbral, y solo se observan los individuos cuyos tiempos de supervivencia superan (o están por debajo de) ese umbral. Existen dos tipos principales de truncamiento: truncamiento a la izquierda y truncamiento a la derecha.

45.3.1 Modelo para Datos Truncados

Para datos truncados a la izquierda, la función de verosimilitud se ajusta para considerar solo los individuos que superan el umbral de truncamiento:

$$L(\beta) = \prod_{i=1}^n \frac{f(t_i | X_i; \beta)}{1 - F(L_i | X_i; \beta)}$$

donde L_i es el umbral de truncamiento para el individuo i .

45.4 Análisis de Competing Risks

En estudios donde pueden ocurrir múltiples tipos de eventos (competing risks), es crucial modelar adecuadamente el riesgo asociado con cada tipo de evento. La probabilidad de ocurrencia de cada evento compite con las probabilidades de ocurrencia de otros eventos.

45.4.1 Modelo de Competing Risks

Para un análisis de competing risks, la función de riesgo se descompone en funciones de riesgo específicas para cada tipo de evento:

$$\lambda(t) = \sum_{j=1}^m \lambda_j(t)$$

donde $\lambda_j(t)$ es la función de riesgo para el evento j .

45.5 Métodos de Imputación

Los métodos de imputación se utilizan para manejar datos faltantes o censurados en estudios de supervivencia. La imputación múltiple es un enfoque común que crea múltiples conjuntos de datos completos imputando valores faltantes varias veces y luego combina los resultados.

45.5.1 Imputación Múltiple

La imputación múltiple para datos de supervivencia se realiza en tres pasos:

1. Imputar los valores faltantes múltiples veces para crear varios conjuntos de datos completos.
2. Analizar cada conjunto de datos completo por separado utilizando métodos de supervivencia estándar.
3. Combinar los resultados de los análisis separados para obtener estimaciones y varianzas combinadas.

45.6 Ejemplo de Análisis con Datos Complicados

Consideremos un estudio con datos censurados por intervalo y competing risks. Ajustamos un modelo para los datos censurados por intervalo y obtenemos los siguientes coeficientes para las covariables edad y tratamiento:

$$\hat{\beta}_{edad} = 0.04, \quad \hat{\beta}_{tratamiento} = -0.8$$

La función de supervivencia ajustada se expresa como:

$$S(t | X) = \exp \left(- \left(\frac{t \exp(0.04 \cdot edad - 0.8 \cdot tratamiento)}{\lambda} \right)^k \right)$$

45.7 Conclusión

El análisis de supervivencia en datos complicados requiere métodos avanzados para manejar censura por intervalo, datos truncados y competing risks. La aplicación de modelos adecuados y métodos de imputación asegura un análisis preciso y completo de estos datos complejos.

Chapter 46

Proyecto Final y Revisión

46.1 Introducción

El proyecto final proporciona una oportunidad para aplicar los conceptos y técnicas aprendidas en el curso de análisis de supervivencia. Este capítulo incluye una guía para desarrollar un proyecto de análisis de supervivencia y una revisión de los conceptos clave.

46.2 Desarrollo del Proyecto

El proyecto final debe incluir los siguientes componentes:

1. Definición del problema: Identificar la pregunta de investigación y los objetivos del análisis de supervivencia.
2. Descripción de los datos: Presentar los datos utilizados, incluyendo las covariables y la estructura de los datos.
3. Análisis exploratorio: Realizar un análisis descriptivo de los datos, incluyendo la censura y la distribución de los tiempos de supervivencia.
4. Ajuste del modelo: Ajustar modelos de supervivencia adecuados (Kaplan-Meier, Cox, AFT) y evaluar su bondad de ajuste.
5. Diagnóstico del modelo: Realizar diagnósticos para evaluar los supuestos del modelo y la influencia de observaciones individuales.
6. Interpretación de resultados: Interpretar los coeficientes del modelo y las curvas de supervivencia ajustadas.
7. Conclusiones: Resumir los hallazgos del análisis y proporcionar recomendaciones basadas en los resultados.

46.3 Revisión de Conceptos Clave

Una revisión de los conceptos clave del análisis de supervivencia incluye:

- **Función de Supervivencia:** Define la probabilidad de sobrevivir más allá de un tiempo específico.
- **Función de Riesgo:** Define la tasa instantánea de ocurrencia del evento.

- **Estimador de Kaplan-Meier:** Proporciona una estimación no paramétrica de la función de supervivencia.
- **Test de Log-rank:** Compara curvas de supervivencia entre diferentes grupos.
- **Modelo de Cox:** Evalúa el efecto de múltiples covariables sobre el tiempo hasta el evento, asumiendo proporcionalidad de riesgos.
- **Modelos AFT:** Modelan el efecto de las covariables multiplicando el tiempo de supervivencia por una constante.
- **Análisis Multivariado:** Considera interacciones y efectos no lineales entre múltiples covariables.
- **Supervivencia en Datos Complicados:** Maneja censura por intervalo, datos truncados y competing risks.

46.4 Ejemplo de Proyecto Final

A continuación se presenta un ejemplo de estructura de un proyecto final de análisis de supervivencia:

46.4.1 Definición del Problema

Analizar el efecto del tratamiento y la edad sobre la supervivencia de pacientes con una enfermedad específica.

46.4.2 Descripción de los Datos

Datos de supervivencia de 100 pacientes, con covariables: edad, sexo y tipo de tratamiento. Los tiempos de supervivencia están censurados a la derecha.

46.4.3 Análisis Exploratorio

Realizar histogramas y curvas de Kaplan-Meier para explorar la distribución de los tiempos de supervivencia y la censura.

46.4.4 Ajuste del Modelo

Ajustar un modelo de Cox y un modelo AFT con las covariables edad y tratamiento.

46.4.5 Diagnóstico del Modelo

Evaluar la proporcionalidad de riesgos y realizar análisis de residuos para identificar observaciones influyentes.

46.4.6 Interpretación de Resultados

Interpretar los coeficientes del modelo y las curvas de supervivencia ajustadas para diferentes niveles de las covariables.

46.4.7 Conclusiones

Resumir los hallazgos y proporcionar recomendaciones sobre el efecto del tratamiento y la edad en la supervivencia de los pacientes.

46.5 Conclusión

El proyecto final es una oportunidad para aplicar los conocimientos adquiridos en un contexto práctico. La revisión de los conceptos clave y la aplicación de técnicas adecuadas de análisis de supervivencia aseguran un análisis riguroso y significativo.

Chapter 47

Documentos Adicionales

47.0.1 Parte I. Introducción a la Bioestadística

47.1 Introducción

47.1.1 Definición de Estadística

- La Estadística es una ciencia formal que estudia la recolección, análisis e interpretación de datos de una muestra representativa, ya sea para ayudar en la toma de decisiones o para explicar condiciones regulares o irregulares de algún fenómeno o estudio aplicado, de ocurrencia en forma aleatoria o condicional.
- Sin embargo, la estadística es más que eso, es decir, es el vehículo que permite llevar a cabo el proceso relacionado con la investigación científica.
- Es transversal a una amplia variedad de disciplinas, desde la física hasta las ciencias sociales, desde las ciencias de la salud hasta el control de calidad. Se usa para la toma de decisiones en áreas de negocios o instituciones gubernamentales.
- Como dijera Huntsberger: *La palabra estadística a menudo nos trae a la mente imágenes de números apilados en grandes arreglos y tablas, de volúmenes de cifras relativas a nacimientos, muertes, impuestos, poblaciones, ingresos, deudas, créditos y así sucesivamente.* Huntsberger tiene razón pues al instante de escuchar esta palabra estas son las imágenes que llegan a nuestra cabeza.
- La Estadística es mucho más que sólo números apilados y gráficas bonitas.
- Es una ciencia con tanta antigüedad como la escritura, y es por sí misma auxiliar de todas las demás ciencias.
- Los mercados, la medicina, la ingeniería, los gobiernos, etc. Se nombran entre los más destacados clientes de ésta.
- La ausencia de ésta conllevaría a un caos generalizado, dejando a los administradores y ejecutivos sin información vital a la hora de tomar decisiones en tiempos de incertidumbre.
- La Estadística que conocemos hoy en día debe gran parte de su realización a los trabajos matemáticos de aquellos hombres que desarrollaron la teoría de las probabilidades, con la cual se adhirió a la Estadística a las ciencias formales.

Definición 32. *La Estadística es la ciencia cuyo objetivo es reunir una información cuantitativa concerniente a individuos, grupos, series de hechos, etc. y deducir de ello gracias al análisis de estos datos unos significados precisos o unas previsiones para el futuro.*

- La estadística, en general, es la ciencia que trata de la recopilación, organización presentación, análisis e interpretación de datos numéricos con el fin de realizar una toma de decisión más efectiva.
- Otros autores la definen como la expresión cuantitativa del conocimiento dispuesta en forma adecuada para el escrutinio y análisis.
- Los estudiantes confunden comúnmente los demás términos asociados con las Estadísticas, una confusión que es conveniente aclarar debido a que esta palabra tiene tres significados:
 - * la palabra estadística, en primer término se usa para referirse a la información estadística;
 - * también se utiliza para referirse al conjunto de técnicas y métodos que se utilizan para analizar la información estadística; y
 - * el término estadístico, en singular y en masculino, se refiere a una medida derivada de una muestra.

47.1.2 Utilidad e Importancia

- Los métodos estadísticos tradicionalmente se utilizan para propósitos descriptivos, para organizar y resumir datos numéricos. La estadística descriptiva, por ejemplo trata de la tabulación de datos, su presentación en forma gráfica o ilustrativa y el cálculo de medidas descriptivas.
- Ahora bien, las técnicas estadísticas se aplican de manera amplia en mercadotecnia, contabilidad, control de calidad y en otras actividades; estudios de consumidores; análisis de resultados en deportes; administradores de instituciones; en la educación; organismos políticos; médicos; y por otras personas que intervienen en la toma de decisiones.

47.1.3 Historia de la Estadística

- Es difícil conocer los orígenes de la Estadística. Desde los comienzos de la civilización han existido formas sencillas de estadística, pues ya se utilizaban representaciones gráficas y otros símbolos en pieles, rocas, palos de madera y paredes de cuevas para contar el número de personas, animales o ciertas cosas.
- Su origen empieza posiblemente en la isla de Cerdeña, donde existen monumentos prehistóricos pertenecientes a los Nuragas, los primeros habitantes de la isla; estos monumentos constan de bloques de basalto superpuestos sin mortero y en cuyas paredes se encontraban grabados toscos signos que han sido interpretados con mucha verosimilitud como muescas que servían para llevar la cuenta del ganado y la caza.
- Los babilonios usaban ya pequeñas tablillas de arcilla para recopilar datos en tablas sobre la producción agrícola y los géneros vendidos o cambiados mediante trueque.
- Otros vestigios pueden ser hallados en el antiguo Egipto, cuyos faraones lograron recopilar, hacia el año 3050 antes de Cristo, prolijos datos relativos a la población y la riqueza del país.
- De acuerdo al historiador griego Heródoto, dicho registro de riqueza y población se hizo con el objetivo de preparar la construcción de las pirámides.
- En el mismo Egipto, Ramsés II hizo un censo de las tierras con el objeto de verificar un nuevo reparto.
- En el antiguo Israel la Biblia da referencias, en el libro de los Números, de los datos estadísticos obtenidos en dos recuentos de la población hebrea. El rey David por otra parte, ordenó a Joab, general del ejército hacer un censo de Israel con la finalidad de conocer el número de la población.

- También los chinos efectuaron censos hace más de cuarenta siglos.
- Los griegos efectuaron censos periódicamente con fines tributarios, sociales (división de tierras) y militares (cálculo de recursos y hombres disponibles).
- La investigación histórica revela que se realizaron 69 censos para calcular los impuestos, determinar los derechos de voto y ponderar la potencia guerrera.
- Fueron los romanos, maestros de la organización política, quienes mejor supieron emplear los recursos de la estadística. Cada cinco años realizaban un censo de la población y sus funcionarios públicos tenían la obligación de anotar nacimientos, defunciones y matrimonios, sin olvidar los recuentos periódicos del ganado y de las riquezas contenidas en las tierras conquistadas.
- Para el nacimiento de Cristo sucedía uno de estos empadronamientos de la población bajo la autoridad del imperio.
- Durante los mil años siguientes a la caída del imperio Romano se realizaron muy pocas operaciones Estadísticas, con la notable excepción de las relaciones de tierras pertenecientes a la Iglesia, compiladas por Pipino el Breve en el 758 y por Carlomagno en el 762 DC.
- Durante el siglo IX se realizaron en Francia algunos censos parciales de siervos.
- En Inglaterra, Guillermo el Conquistador recopiló el Domesday Book o libro del Gran Catastro para el año 1086, un documento de la propiedad, extensión y valor de las tierras de Inglaterra. Esa obra fue el primer compendio estadístico de Inglaterra.
- Aunque Carlomagno, en Francia; y Guillermo el Conquistador, en Inglaterra, trataron de revivir la técnica romana, los métodos estadísticos permanecieron casi olvidados durante la Edad Media.
- Durante los siglos XV, XVI, y XVII, hombres como Leonardo de Vinci, Nicolás Copérnico, Galileo, Neper, William Harvey, Sir Francis Bacon y René Descartes, hicieron grandes operaciones al método científico, de tal forma que cuando se crearon los Estados Nacionales y surgió como fuerza el comercio internacional existía ya un método capaz de aplicarse a los datos económicos.
- Para el año 1532 empezaron a registrarse en Inglaterra las defunciones debido al temor que Enrique VII tenía por la peste. Más o menos por la misma época, en Francia la ley exigió a los clérigos registrar los bautismos, fallecimientos y matrimonios.
- Durante un brote de peste que apareció a fines de la década de 1500, el gobierno inglés comenzó a publicar estadística semanales de los decesos. Esa costumbre continuó muchos años, y en 1632 estos Bills of Mortality (Cuentas de Mortalidad) contenían los nacimientos y fallecimientos por sexo.
- En 1662, el capitán John Graunt usó documentos que abarcaban treinta años y efectuó predicciones sobre el número de personas que morirían de varias enfermedades y sobre las proporciones de nacimientos de varones y mujeres que cabría esperar.
- El trabajo de Graunt, condensado en su obra *Natural and Political Observations...Made upon the Bills of Mortality*, fue un esfuerzo innovador en el análisis estadístico. Por el año 1540 el alemán Sebastián Muster realizó una compilación estadística de los recursos nacionales, comprensiva de datos sobre organización política, instrucciones sociales, comercio y poderío militar.
- Durante el siglo XVII aportó indicaciones más concretas de métodos de observación y análisis cuantitativo y amplió los campos de la inferencia y la teoría Estadística.
- Los eruditos del siglo XVII demostraron especial interés por la Estadística Demográfica como resultado de la especulación sobre si la población aumentaba, decrecía o permanecía estática.
- En los tiempos modernos tales métodos fueron resucitados por algunos reyes que necesitaban conocer las riquezas monetarias y el potencial humano de sus respectivos países.

- El primer empleo de los datos estadísticos para fines ajenos a la política tuvo lugar en 1691 y estuvo a cargo de Gaspar Neumann, un profesor alemán que vivía en Breslau.
- Este investigador se propuso destruir la antigua creencia popular de que en los años terminados en siete moría más gente que en los restantes, y para lograrlo hurgó pacientemente en los archivos parroquiales de la ciudad.
- Después de revisar miles de partidas de defunción pudo demostrar que en tales años no fallecían más personas que en los demás.
- Los procedimientos de Neumann fueron conocidos por el astrónomo inglés Halley, descubridor del cometa que lleva su nombre, quien los aplicó al estudio de la vida humana.
- Sus cálculos sirvieron de base para las tablas de mortalidad que hoy utilizan todas las compañías de seguros.
- Durante el siglo XVII y principios del XVIII, matemáticos como Bernoulli, Francis Maseres, Lagrange y Laplace desarrollaron la teoría de probabilidades.
- No obstante durante cierto tiempo, la teoría de las probabilidades limitó su aplicación a los juegos de azar y hasta el siglo XVIII no comenzó a aplicarse a los grandes problemas científicos.
- Godofredo Achenwall, profesor de la Universidad de Gotinga, acuñó en 1760 la palabra estadística, que extrajo del término italiano statista (estadista).
- Creía, y con sobrada razón, que los datos de la nueva ciencia serían el aliado más eficaz del gobernante consciente.
- La raíz remota de la palabra se halla, por otra parte, en el término latino status, que significa estado o situación; Esta etimología aumenta el valor intrínseco de la palabra, por cuanto la estadística revela el sentido cuantitativo de las más variadas situaciones.
- Jacques Quételet es quien aplica las Estadísticas a las ciencias sociales. Este interpretó la teoría de la probabilidad para su uso en las ciencias sociales y resolver la aplicación del principio de promedios y de la variabilidad a los fenómenos sociales.
- Quételet fue el primero en realizar la aplicación práctica de todo el método Estadístico, entonces conocido, a las diversas ramas de la ciencia.
- Entretanto, en el período del 1800 al 1820 se desarrollaron dos conceptos matemáticos fundamentales para la teoría Estadística; la teoría de los errores de observación, aportada por Laplace y Gauss; y la teoría de los mínimos cuadrados desarrollada por Laplace, Gauss y Legendre.
- A finales del siglo XIX, Sir Francis Gaston ideó el método conocido por Correlación, que tenía por objeto medir la influencia relativa de los factores sobre las variables.
- De aquí partió el desarrollo del coeficiente de correlación creado por Karl Pearson y otros cultivadores de la ciencia biométrica como J. Pease Norton, R. H. Hooker y G. Udny Yule, que efectuaron amplios estudios sobre la medida de las relaciones.
- Los progresos más recientes en el campo de la Estadística se refieren al ulterior desarrollo del cálculo de probabilidades, particularmente en la rama denominada indeterminismo o relatividad, se ha demostrado que el determinismo fue reconocido en la Física como resultado de las investigaciones atómicas y que este principio se juzga aplicable tanto a las ciencias sociales como a las físicas.

47.1.4 Etapas de Desarrollo de la Estadística

La historia de la estadística está resumida en tres grandes etapas o fases.

- Fase 1: **Los Censos:** Desde el momento en que se constituye una autoridad política, la idea de inventariar de una forma más o menos regular la población y las riquezas existentes en el territorio está ligada a la conciencia de soberanía y a los primeros esfuerzos administrativos.
- Fase 2: **De la Descripción de los Conjuntos a la Aritmética Política:** Las ideas mercantilistas extrañan una intensificación de este tipo de investigación. Colbert multiplica las encuestas sobre artículos manufacturados, el comercio y la población: los intendentes del Reino envían a París sus memorias. Vauban, más conocido por sus fortificaciones o su Dime Royale, que es la primera propuesta de un impuesto sobre los ingresos, se señala como el verdadero precursor de los sondeos. Más tarde, Bufón se preocupa de esos problemas antes de dedicarse a la historia natural. La escuela inglesa proporciona un nuevo progreso al superar la fase puramente descriptiva.
- Fase 2: Sus tres principales representantes son Graunt, Petty y Halley. El penúltimo es autor de la famosa Aritmética Política. Chaptal, ministro del interior francés, publica en 1801 el primer censo general de población, desarrolla los estudios industriales, de las producciones y los cambios, haciéndose sistemáticos durante las dos terceras partes del siglo XIX.
- Fase 3: **Estadística y Cálculo de Probabilidades:** El cálculo de probabilidades se incorpora rápidamente como un instrumento de análisis extremadamente poderoso para el estudio de los fenómenos económicos y sociales y en general para el estudio de fenómenos cuyas causas son demasiado complejas para conocerlos totalmente y hacer posible su análisis.

47.1.5 División de la Estadística

La Estadística para su mejor estudio se ha dividido en dos grandes ramas: **la Estadística Descriptiva y la Estadística Inferencial.**

- **Descriptiva:** consiste sobre todo en la presentación de datos en forma de tablas y gráficas. Esta comprende cualquier actividad relacionada con los datos y está diseñada para resumir o describir los mismos sin factores pertinentes adicionales; esto es, sin intentar inferir nada que vaya más allá de los datos, como tales.
- **Inferencial:** se deriva de muestras, de observaciones hechas sólo acerca de una parte de un conjunto numeroso de elementos y esto implica que su análisis requiere de generalizaciones que van más allá de los datos. Como consecuencia, la característica más importante del reciente crecimiento de la estadística ha sido un cambio en el énfasis de los métodos que describen a métodos que sirven para hacer generalizaciones. La Estadística Inferencial investiga o analiza una población partiendo de una muestra tomada.

47.1.6 Estadística Inferencial

Los métodos básicos de la estadística inferencial son la estimación y el contraste de hipótesis, que juegan un papel fundamental en la investigación. Por tanto, algunos de los objetivos que se persiguen son:

- Calcular los parámetros de la distribución de medias o proporciones muestrales de tamaño n , extraídas de una población de media y varianza conocidas.
- Estimar la media o la proporción de una población a partir de la media o proporción muestral.
- Utilizar distintos tamaños muestrales para controlar la confianza y el error admitido.
- Contrastar los resultados obtenidos a partir de muestras.
- Visualizar gráficamente, mediante las respectivas curvas normales, las estimaciones realizadas.

- En definitiva, la idea es, a partir de una población se extrae una muestra por algunos de los métodos existentes, con la que se generan datos numéricos que se van a utilizar para generar estadísticos con los que realizar estimaciones o contrastes poblacionales.
- Existen dos formas de estimar parámetros: la *estimación puntual* y la *estimación por intervalo de confianza*. En la primera se busca, con base en los datos muestrales, un único valor estimado para el parámetro. Para la segunda, se determina un intervalo dentro del cual se encuentra el valor del parámetro, con una probabilidad determinada.
- Si el objetivo del tratamiento estadístico inferencial, es efectuar generalizaciones acerca de la estructura, composición o comportamiento de las poblaciones no observadas, a partir de una parte de la población, será necesario que la parcela de población examinada sea representativa del total.
- Por ello, la selección de la muestra requiere unos requisitos que lo garanticen, debe ser representativa y aleatoria.
- Además, la cantidad de elementos que integran la muestra (el tamaño de la muestra) depende de múltiples factores, como el dinero y el tiempo disponibles para el estudio, la importancia del tema analizado, la confiabilidad que se espera de los resultados, las características propias del fenómeno analizado, etcétera.
- Así, a partir de la muestra seleccionada se realizan algunos cálculos y se estima el valor de los parámetros de la población tales como la media, la varianza, la desviación estándar, o la forma de la distribución, etc.

47.1.7 Método Estadístico

El conjunto de los métodos que se utilizan para medir las características de la información, para resumir los valores individuales, y para analizar los datos a fin de extraerles el máximo de información, es lo que se llama *métodos estadísticos*. Los métodos de análisis para la información cuantitativa se pueden dividir en los siguientes seis pasos:

1. Definición del problema.
 2. Recopilación de la información existente.
 3. Obtención de información original.
 4. Clasificación.
 5. Presentación.
 6. Análisis.
- El centro de gravedad de la metodología estadística se empieza a desplazar técnicas de computación intensiva aplicadas a grandes masas de datos, y se empieza a considerar el método estadístico como un proceso iterativo de búsqueda del modelo ideal.
 - Las aplicaciones en este periodo de la Estadística a la Economía conducen a una disciplina con contenido propio: la Econometría. La investigación estadística en problemas militares durante la segunda guerra mundial y los nuevos métodos de programación matemática, dan lugar a la Investigación Operativa

47.1.8 Errores Estadísticos Comunes

Al momento de recopilar los datos que serán procesados se es susceptible de cometer errores así como durante los cálculos de los mismos. No obstante, hay otros errores que no tienen nada que ver con la digitación y que no son tan fácilmente identificables. Algunos de éstos errores son:

- **Sesgo:** Es imposible ser completamente objetivo o no tener ideas preconcebidas antes de comenzar a estudiar un problema, y existen muchas maneras en que una perspectiva o estado mental pueda influir en la recopilación y en el análisis de la información. En estos casos se dice que hay un sesgo cuando el individuo da mayor peso a los datos que apoyan su opinión que a aquellos que la contradicen. Un caso extremo de sesgo sería la situación donde primero se toma una decisión y después se utiliza el análisis estadístico para justificar la decisión ya tomada.
- **Datos No Comparables:** el establecer comparaciones es una de las partes más importantes del análisis estadístico, pero es extremadamente importante que tales comparaciones se hagan entre datos que sean comparables.
- **Proyección descuidada de tendencias:** la proyección simplista de tendencias pasadas hacia el futuro es uno de los errores que más ha desacreditado el uso del análisis estadístico.
- **Muestreo Incorrecto:** en la mayoría de los estudios sucede que el volumen de información disponible es tan inmenso que se hace necesario estudiar muestras, para derivar conclusiones acerca de la población a que pertenece la muestra. Si la muestra se selecciona correctamente, tendrá básicamente las mismas propiedades que la población de la cual fue extraída; pero si el muestreo se realiza incorrectamente, entonces puede suceder que los resultados no signifiquen nada

En resumen se puede decir que la Estadística es un conjunto de procedimientos para reunir, clasificar, codificar, procesar, analizar y resumir información numérica adquirida sistemáticamente (Ritchey, 2002). Permite hacer inferencias a partir de una muestra para extrapolarlas a una población. Aunque normalmente se asocia a muchos cálculos y operaciones aritméticas, y aunque las matemáticas están involucradas, en su mayor parte sus fundamentos y uso apropiado pueden dominarse sin hacer referencia a habilidades matemáticas avanzadas.

De hecho se trata de una forma de ver la realidad basada en el análisis cuidadoso de los hechos (Ritchey, 2002). Es necesaria sin embargo la sistematización para reducir el efecto que las emociones y las experiencias individuales puedan tener al interpretar esa realidad.

De esta manera la estadística se relaciona con el método científico complementándolo como herramienta de análisis y, aunque la investigación científica no requiere necesariamente de la estadística, ésta valida muchos de los resultados cuantitativos derivados de la investigación.

La obtención del conocimiento debe hacerse de manera sistemática por lo que deben planearse todos los pasos que llevan desde el planteamiento de un problema, pasando por la elaboración de hipótesis y la manera en que van a ser probadas; la selección de sujetos (muestreo), los escenarios, los instrumentos que se utilizarán para obtener los datos, definir el procedimiento que se seguirá para esto último, los controles que se deben hacer para asegurar que las intervenciones son las causas más probables de los cambios esperados (diseño);

El tratamiento de los datos de la investigación científica tiene varias etapas:

- En la etapa de recolección de datos del método científico, se define a la población de interés y se selecciona una muestra o conjunto de personas representativas de la misma, se realizan experimentos o se emplean instrumentos ya existentes o de nueva creación, para medir los atributos de interés necesarios para responder a las preguntas de investigación. Durante lo que es llamado trabajo de campo se obtienen los datos en crudo, es decir las respuestas directas de los sujetos uno por uno, se codifican (se les asignan valores a las respuestas), se capturan y se verifican para ser utilizados en las siguientes etapas.

- En la etapa de recuento, se organizan y ordenan los datos obtenidos de la muestra. Esta será descrita en la siguiente etapa utilizando la estadística descriptiva, todas las investigaciones utilizan estadística descriptiva, para conocer de manera organizada y resumida las características de la muestra.
- En la etapa de análisis se utilizan las pruebas estadísticas (estadística inferencial) y en la interpretación se acepta o rechaza la hipótesis nula.
- En investigación, el fenómeno en estudio puede ser cualitativo que implicaría comprenderlo y explicarlo, o cuantitativo para compararlo y hacer inferencias. Se puede decir que si se hace análisis se usan métodos cuantitativos y si se hace descripción se usan métodos cualitativos.

Medición Para poder emplear el método estadístico en un estudio es necesario medir las variables.

- Medir: es asignar valores a las propiedades de los objetos bajo ciertas reglas, esas reglas son los niveles de medición
- Cuantificar: es asignar valores a algo tomando un patrón de referencia. Por ejemplo, cuantificar es ver cuántos hombres y cuántas mujeres hay.
- Variable: es una característica o propiedad que asume diferentes valores dentro de una población de interés y cuya variación es susceptible de medirse. Las variables pueden clasificarse de acuerdo al tipo de valores que puede tomar como:
 - * Discretas o categóricas.- en las que los valores se relacionan a nombres, etiquetas o categorías, no existe un significado numérico directo
 - * Continuas.- los valores tienen un correlato numérico directo, son continuos y susceptibles de fraccionarse y de poder utilizarse en operaciones aritméticas De acuerdo a la cantidad de valores
 - * Dicotómica.- sólo tienen dos valores posibles, la característica está ausente o presente
 - * Policotómica.- pueden tomar tres valores o más, pueden tomarse matices diferentes, en grados, jerarquías o magnitudes continuas.
- En cuanto a una clasificación estadística
 - * Aleatoria.- Aquella en la cual desconocemos el valor porque fluctúa de acuerdo a un evento debido al azar
 - * Determinística.- Aquella variable de la que se conoce el valor
 - * Independiente.- aquellas variables que son manipuladas por el investigador. Define los grupos
 - * Dependiente.- son mediciones que ocurren durante el experimento o tratamiento (resultado de la independiente), es la que se mide y compara entre los grupos

Niveles de Medición

- Nominal Las propiedades de la medición nominal son:
 - * Exhaustiva: implica a todas las opciones
 - * A los sujetos se les asignan categorías, por lo que son mutuamente excluyentes. Es decir, la variable está presente o no; tiene o no una característica
- Ordinal Las propiedades de la medición ordinal son:
 - * El nivel ordinal posee transitividad, por lo que se tiene la capacidad de identificar que es mejor o mayor que otra, en ese sentido se pueden establecer jerarquías
 - * Las distancias entre un valor y otro no son iguales.
- Intervalo

- * El nivel de medición intervalar requiere distancias iguales entre cada valor. Por lo general utiliza datos cuantitativos. Por ejemplo: temperatura, atributos psicológicos (CI, nivel de autoestima, pruebas de conocimientos, etc.)
 - * Las unidades de calificación son equivalentes en todos los puntos de la escala. Una escala de intervalos implica: clasificación, magnitud y unidades de tamaños iguales (Brown, 2000).
 - * Se pueden hacer operaciones aritméticas
 - * Cuando se le pide al sujeto que califique una situación del 0 al 10 puede tomarse como un nivel de medición de intervalo, siempre y cuando se incluya el 0.
- Razón
- * La escala empieza a partir del 0 absoluto, por lo tanto incluye sólo los números por su valor en sí, por lo que no pueden existir los números con signo negativo. Por ejemplo: Peso corporal en kg., edad en años, estatura en cm.
 - * Convencionalmente los datos que son de nivel absoluto o de razón son manejados como los datos intervalares.

47.1.9 Términos comunes utilizados en Estadística

- **Variable:** Consideraciones que una variable son una característica o fenómeno que puede tomar distintos valores.
- **Dato:** Mediciones o cualidades que han sido recopiladas como resultado de observaciones.
- **Población:** Se considera el área de la cual son extraídos los datos. Es decir, es el conjunto de elementos o individuos que poseen una característica común y medible acerca de lo cual se desea información. Es también llamado Universo.
- **Muestra:** Es un subconjunto de la población, seleccionado de acuerdo a una regla o algún plan de muestreo.
- **Censo:** Recopilación de todos los datos (de interés para la investigación) de la población.
- **Estadística:** Es una función o fórmula que depende de los datos de la muestra (es variable).
- **Parámetro:** Característica medible de la población. Es un resumen numérico de alguna variable observada de la población. Los parámetros normales que se estudian son: *La media poblacional, la media poblacional, Proporción.*
- **Estimador:** Un estimador de un parámetro, es un estadístico que se emplea para conocer el parámetro desconocido.
- **Estadístico:** Es una función de los valores de la muestra. Es una variable aleatoria, cuyos valores dependen de la muestra seleccionada. Su distribución de probabilidad, se conoce como *Distribución muestral del estadístico.*
- **Estimación:** Este término indica que a partir de lo observado en una muestra (un resumen estadístico con las medidas que conocemos de Descriptiva) se extrapola o generaliza dicho resultado muestral a la población total, de modo que lo estimado es el valor generalizado a la población. Consiste en la búsqueda del valor de los parámetros poblacionales objeto de estudio. Puede ser puntual o por intervalo de confianza:
 - * *Puntual:* cuando buscamos un valor concreto. Un estimador de un parámetro poblacional es una función de los datos muestrales. En pocas palabras, es una fórmula que depende de los valores obtenidos de una muestra, para realizar estimaciones. Lo que se pretende obtener es el valor exacto de un parámetro.

Las propiedades deseables de un estimador son los siguientes:

- Insesgado: Diremos que un estimador de un parámetro es insesgado si su esperanza coincide con el verdadero valor del parámetro. En el caso de que no coincidan, diremos que el estimador es sesgado.
- Eficiencia: Dados dos estimadores para un mismo parámetro, se dice que uno es más eficiente que el otro si tiene menor varianza.
- Suficiencia: Se dice que un estimador de un parámetro es suficiente cuando para su cálculo utiliza toda la información de la muestra.
- Consistencia: Decimos que un estimador de un parámetro es consistente si la distribución del estimador tiende a concentrarse en un cierto punto cuando el tamaño de la muestra tiende a infinito.

Demostrar que un cierto estimador cumple estas propiedades puede ser complicado en determinadas ocasiones. Existen varios métodos que nos van a permitir obtener los estimadores puntuales. Los más importantes son:

- Método de Momentos: se basa en que los momentos poblacionales y se estiman mediante los momentos muestrales. Suelen dar estimadores consistentes.
- Método de Mínimos Cuadrados: consiste en obtener un estimador que hace mínima una determinada función.
- Método de Máxima Verosimilitud: consiste en tomar como parámetro poblacional el valor de la muestra que sea más probable, es decir, que tenga mayor probabilidad. Se suelen obtener estimadores consistentes y eficientes. Es el más utilizado.
- *Intervalo de confianza*: cuando determinamos un intervalo, dentro del cual se supone que va a estar el valor del parámetro que se busca con una cierta probabilidad. El intervalo de confianza está determinado por dos valores dentro de los cuales afirmamos que está el verdadero parámetro con cierta probabilidad. Son unos límites o margen de variabilidad que damos al valor estimado, para poder afirmar, bajo un criterio de probabilidad, que el verdadero valor no los rebasará. Este intervalo contiene al parámetro estimado con una determinada certeza o nivel de confianza.

En la estimación por intervalos se usan los siguientes conceptos:

- Variabilidad del parámetro: Si no se conoce, puede obtenerse una aproximación en los datos o en un estudio piloto. También hay métodos para calcular el tamaño de la muestra que prescindan de este aspecto. Habitualmente se usa como medida de esta variabilidad la desviación típica poblacional.
- Error de la estimación: Es una medida de su precisión que se corresponde con la amplitud del intervalo de confianza. Cuanta más precisión se desee en la estimación de un parámetro, más estrecho deberá ser el intervalo de confianza y, por tanto, menor el error, y más sujetos deberán incluirse en la muestra estudiada.
- Nivel de confianza: Es la probabilidad de que el verdadero valor del parámetro estimado en la población se sitúe en el intervalo de confianza obtenido. El nivel de confianza se denota por $1 - \alpha$
- *p-value*: También llamado nivel de significación. Es la probabilidad (en tanto por uno) de fallar en nuestra estimación, esto es, la diferencia entre la certeza (1) y el nivel de confianza $1 - \alpha$.
- Valor crítico: Se representa por $Z_{\alpha/2}$. Es el valor de la abscisa en una determinada distribución que deja a su derecha un área igual a $1/2$, siendo $1 - \alpha$ el nivel de confianza. Normalmente los valores críticos están tabulados o pueden calcularse en función de la distribución de la población.

Para un tamaño fijo de la muestra, los conceptos de error y nivel de confianza van relacionados. Si admitimos un error mayor, esto es, aumentamos el tamaño del intervalo de confianza, tenemos también una mayor probabilidad de éxito en nuestra estimación, es decir, un mayor nivel de confianza. Por tanto, un aspecto que debe de tenerse en cuenta es el tamaño muestral, ya que para disminuir el error que se comente habrá que aumentar el tamaño muestral. Esto se resolverá, para un intervalo de confianza cualquiera, despejando el tamaño de la muestra en cualquiera de las formulas de los intervalos de confianza que veremos a continuación, a partir del error máximo permitido. Los intervalos de confianza pueden ser unilaterales o bilaterales:

- **Contraste de Hipótesis:** Consiste en determinar si es aceptable, partiendo de datos muestrales, que la característica o el parámetro poblacional estudiado tome un determinado valor o esté dentro de unos determinados valores.
- **Nivel de Confianza:** Indica la proporción de veces que acertaríamos al afirmar que el parámetro está dentro del intervalo al seleccionar muchas muestras.

47.1.10 Muestreo:

Una muestra es representativa en la medida que es imagen de la población.

En general, podemos decir que el tamaño de una muestra dependerá principalmente de: *Nivel de precisión deseado, Recursos disponibles, Tiempo involucrado en la investigación..* Además el plan de muestreo debe considerar *La población, Parámetros a medir.*

Existe una gran cantidad de tipos de muestreo. En la práctica los más utilizados son los siguientes:

- **MUESTREO ALEATORIO SIMPLE:** Es un método de selección de n unidades extraídas de N , de tal manera que cada una de las posibles muestras tiene la misma probabilidad de ser escogida. (En la práctica, se enumeran las unidades de 1 a N , y a continuación se seleccionan n números aleatorios entre 1 y N , ya sea de tablas o de alguna urna con fichas numeradas).
- **MUESTREO ESTRATIFICADO ALEATORIO:** Se usa cuando la población está agrupada en pocos estratos, cada uno de ellos son muchas entidades. Este muestreo consiste en sacar una muestra aleatoria simple de cada uno de los estratos. (Generalmente, de tamaño proporcional al estrato).
- **MUESTREO SISTEMÁTICO :** Se utiliza cuando las unidades de la población están de alguna manera totalmente ordenadas. Para seleccionar una muestra de n unidades, se divide la población en n subpoblaciones de tamaño $K = N/n$ y se toma al azar una unidad de la K primeras y de ahí en adelante cada K -ésima unidad.
- **MUESTREO POR CONGLOMERADO:** Se emplea cuando la población está dividida en grupos o conglomerados pequeños. Consiste en obtener una muestra aleatoria simple de conglomerados y luego CENSAR cada uno de éstos.
- **MUESTREO EN DOS ETAPAS (Bietápico):** En este caso la muestra se toma en dos pasos:
 - * Seleccionar una muestra de unidades primarias, y Seleccionar una muestra de elementos a partir de cada unidad primaria escogida.
 - * Observación: En la realidad es posible encontrarse con situaciones en las cuales no es posible aplicar libremente un tipo de muestreo, incluso estaremos obligados a mezclarlas en ocasiones.

47.1.11 Variables

Las variables se pueden clasificar en dos grandes grupos.

- **Variables categóricas:** Son aquellas que pueden ser representadas a través de símbolos, letras, palabras, etc. Los valores que toman se denominan categorías, y los elementos que pertenecen a estas categorías, se consideran idénticos respecto a la característica que se está midiendo. Las variables categóricas se dividen en dos tipos: Ordinal y Nominal.
 - * **Las Ordinales,** son aquellas en que las categorías tienen un orden implícito. Admiten grados de calidad, es decir, existe una relación total entre las categorías.
 - * **Las nominales,** son aquellas donde no existe una relación de orden.
- **Variables Numéricas:** Son aquellas que pueden tomar valores numéricos exclusivamente (mediciones). Se dividen en dos tipos: Discretas y continuas.
 - * **Discretas:** son aquellas que toman sus valores en un conjunto finito o infinito numerable.
 - * **Continuas:** Son aquellas que toman sus valores en un subconjunto de los números reales, es decir en un intervalo. En general para las variables continuas el hombre ha debido inventar una medida para poder establecer una medición de ellas.

47.1.12 Malos Usos de la Estadística

El propósito de esta sección es solamente indicar los malos usos comunes de datos estadísticos, sin incluir el uso de métodos estadísticos complicados. Un estudiante debería estar alerta en relación con estos malos usos y debería hacer un gran esfuerzo para evitarlos a fin de ser un verdadero estadístico.

Datos estadísticos inadecuados Los datos estadísticos son usados como la materia prima para un estudio estadístico. Cuando los datos son inadecuados, la conclusión extraída del estudio de los datos se vuelve obviamente inválida. Por ejemplo, supongamos que deseamos encontrar el ingreso familiar típico del año pasado en la ciudad Y de 50,000 familias y tenemos una muestra consistente del ingreso de solamente tres familias: 1 millón, 2 millones y no ingreso. Si sumamos el ingreso de las tres familias y dividimos el total por 3, obtenemos un promedio de 1 millón.

Entonces, extraemos una conclusión basada en la muestra de que el ingreso familiar promedio durante el año pasado en la ciudad fue de 1 millón.

Es obvio que la conclusión es falsa, puesto que las cifras son extremas y el tamaño de la muestra es demasiado pequeño; por lo tanto la muestra no es representativa.

Hay muchas otras clases de datos inadecuados. Por ejemplo, algunos datos son respuestas inexactas de una encuesta, porque las preguntas usadas en la misma son vagas o engañosas, algunos datos son toscas estimaciones porque no hay disponibles datos exactos o es demasiado costosa su obtención, y algunos datos son irrelevantes en un problema dado, porque el estudio estadístico no está bien planeado.

47.1.13 Un sesgo del usuario

Sesgo significa que un usuario dé los datos perjudicialmente de más énfasis a los hechos, los cuales son empleados para mantener su predeterminada posición u opinión.

Los estadísticos son frecuentemente degradados por lemas tales como :Hay tres clases de mentiras: mentiras, mentiras reprobables y estadística, y Las cifras no mienten, pero los mentirosos piensan.

Hay dos clases de sesgos: conscientes e inconscientes. Ambos son comunes en el análisis estadístico. Hay numerosos ejemplos de sesgos conscientes.

Un anunciante frecuentemente usa la estadística para probar que su producto es muy superior al producto de su competidor.

Un político prefiere usar la estadística para sostener su punto de vista.

Gerentes y líderes de trabajadores pueden simultáneamente situar sus respectivas cifras estadísticas sobre la misma tabla de trato para mostrar que sus rechazos o peticiones son justificadas.

Es casi imposible que un sesgo inconsciente esté completamente ausente en un trabajo estadístico.

En lo que respecta al ser humano, es difícil obtener una actitud completamente objetiva al abordar un problema, aun cuando un científico debería tener una mente abierta.

Un estadístico debería estar enterado del hecho de que su interpretación de los resultados del análisis estadístico está influenciado por su propia experiencia, conocimiento y antecedentes con relación al problema dado.

47.1.14 Supuestos falsos

Es muy frecuente que un análisis estadístico contemple supuestos. Un investigador debe ser muy cuidadoso en este hecho, para evitar que éstos sean falsos. Los supuestos falsos pueden ser originados por:

- Quien usa los datos
- Quien está tratando de confundir (con intencionalidad)
- Ignorancia
- Descuido.

TÉRMINOS COMUNES UTILIZADOS EN ESTADÍSTICA

- Variable: Consideraciones que una variable son una característica o fenómeno que puede tomar distintos valores.
- Dato: Mediciones o cualidades que han sido recopiladas como resultado de observaciones.
- Población: Se considera el área de la cual son extraídos los datos. Es decir, es el conjunto de elementos o individuos que poseen una característica común y medible acerca de lo cual se desea información. Es también llamado Universo.
- Muestra: Es un subconjunto de la población, seleccionado de acuerdo a una regla o algún plan de muestreo.
- Censo: Recopilación de todos los datos (de interés para la investigación) de la población.
- Estadística: Es una función o fórmula que depende de los datos de la muestra (es variable).
- Parámetro Característica medible de la población.

Ejemplo: La universidad está interesada en determinar el ingreso de las familias de sus alumnos.

- Variable: Ingreso per cápita de las familias.
- Dato: Ingreso per cápita de la familia de un alumno específico.
- Población: Las familias de todos los alumnos de la universidad.
- Estadística: Ingreso per cápita promedio de las familias seleccionadas en la muestra.
- Parámetro: Ingreso per cápita promedio de la población.

MUESTREO Una muestra es representativa en la medida que es imagen de la población. En general, podemos decir que el tamaño de una muestra dependerá principalmente de:

- Nivel de precisión deseado.
- Recursos disponibles.
- Tiempo involucrado en la investigación.

Además el plan de muestreo debe considerar

- La población
- Parámetros a medir.

Existe una gran cantidad de tipos de muestreo. En la práctica los más utilizados son los siguientes:

- **MUESTREO ALEATORIO SIMPLE:** Es un método de selección de n unidades extraídas de N , de tal manera que cada una de las posibles muestras tiene la misma probabilidad de ser escogida. (En la práctica, se enumeran las unidades de 1 a N , y a continuación se seleccionan n números aleatorios entre 1 y N , ya sea de tablas o de alguna urna con fichas numeradas).
- **MUESTREO ESTRATIFICADO ALEATORIO:** Se usa cuando la población está agrupada en pocos estratos, cada uno de ellos son muchas entidades. Este muestreo consiste en sacar una muestra aleatoria simple de cada uno de los estratos. (Generalmente, de tamaño proporcional al estrato).
- **MUESTREO SISTEMÁTICO :** Se utiliza cuando las unidades de la población están de alguna manera totalmente ordenadas. Para seleccionar una muestra de n unidades, se divide la población en n subpoblaciones de tamaño $K = N/n$ y se toma al azar una unidad de la K primeras y de ahí en adelante cada K – ésima unidad, es decir, siendo n o la primera unidad seleccionada de la sub-población $(1, 2, \dots, K)$. $\{n, n + K, n + 2K, \dots, n + (n - 1)K\}$
- **MUESTREO POR CONGLOMERADO** Se emplea cuando la población está dividida en grupos o conglomerados pequeños. Consiste en obtener una muestra aleatoria simple de conglomerados y luego CENSAR cada uno de éstos.
- **MUESTREO EN DOS ETAPAS (Bietápico)** En este caso la muestra se toma en dos pasos:
 - * Seleccionar una muestra de unidades primarias, y
 - * Seleccionar una muestra de elementos a partir de cada unidad primaria escogida.

Observación: En la realidad es posible encontrarse con situaciones en las cuales no es posible aplicar libremente un tipo de muestreo, incluso estaremos obligados a mezclarlas en ocasiones. En general la Estadística está encargada de llevar a cabo el siguiente esquema:

- Recopilar
- Organizar
- Presentar
- Analizar

Tipos de variables Las variables se pueden clasificar en dos grandes grupos.

- **categorías:** Son aquellas que pueden ser representadas a través de símbolos, letras, palabras, etc. Los valores que toman se denominan categorías, y los elementos que pertenecen a estas categorías, se consideran idénticos respecto a la característica que se está midiendo.

Las variables categóricas se dividen en dos tipos: Ordinal y Nominal.

- * Las Ordinales, son aquellas en que las categorías tienen un orden implícito. Admiten grados de calidad, es decir, existe una relación total entre las categorías. A pesar de que esta variable admite grados de calidad, no es posible cuantificar la diferencia.

- * Las nominales, son aquellas donde no existe una relación de orden.
- Variables numéricas. Son aquellas que pueden tomar valores numéricos exclusivamente (mediciones). dividen en dos tipos. Discretas y continuas.
 - * Discretas: son aquellas que toman sus valores en un conjunto finito o infinito numerable.
 - * Continuas: Son aquellas que toman sus valores en un subconjunto de los números reales, es decir en un intervalo.

Observación: En general para las variables continuas el hombre ha debido inventar una medida para poder establecer una medición de ellas: Ejemplo: El metro, la hora.

Los métodos básicos de la estadística inferencial son la estimación y el contraste de hipótesis, que juegan un papel fundamental en la investigación. Por tanto, algunos de los objetivos que se persiguen en este tema son: Inferencia, estimación y contraste de hipótesis

- Calcular los parámetros de la distribución de medias o proporciones muestrales de tamaño n , extraídas de una población de media y varianza conocidas.
- Estimar la media o la proporción de una población a partir de la media o proporción muestral.
- Utilizar distintos tamaños muestrales para controlar la confianza y el error admitido.
- Contrastar los resultados obtenidos a partir de muestras.
- Visualizar gráficamente, mediante las respectivas curvas normales, las estimaciones realizadas.

En la mayoría de las investigaciones resulta imposible estudiar a todos y cada uno de los individuos de la población ya sea por el coste que supondría, o por la imposibilidad de acceder a ello. Mediante la técnica inferencial obtendremos conclusiones para una población no observada en su totalidad, a partir de estimaciones o resúmenes numéricos efectuados sobre la base informativa extraída de una muestra de dicha población.

- Por tanto, el esquema que se sigue es, a partir de una población se extrae una muestra por algunos de los métodos existentes, con la que se generan datos numéricos que se van a utilizar para generar estadísticos con los que realizar estimaciones o contrastes poblacionales.
- Existen dos formas de estimar parámetros: la estimación puntual y la estimación por intervalo de confianza. En la primera se busca, con base en los datos muestrales, un único valor estimado para el parámetro. Para la segunda, se determina un intervalo dentro del cual se encuentra el valor del parámetro, con una probabilidad determinada.
- Si el objetivo del tratamiento estadístico inferencial, es efectuar generalizaciones acerca de la estructura, composición o comportamiento de las poblaciones no observadas, a partir de una parte de la población, será necesario que la parcela de población examinada sea representativa del total. Por ello, la selección de la muestra requiere unos requisitos que lo garanticen, debe ser representativa y aleatoria.

Además, la cantidad de elementos que integran la muestra (el tamaño de la muestra) depende de múltiples factores, como el dinero y el tiempo disponibles para el estudio, la importancia del tema analizado, la confiabilidad que se espera de los resultados, las características propias del fenómeno analizado, etcétera. Así, a partir de la muestra seleccionada se realizan algunos cálculos y se estima el valor de los parámetros de la población tales como la media, la varianza, la desviación estándar, o la forma de la distribución, etc.

Conceptos básicos

- POBLACIÓN: Conjunto de elementos sobre los que se observa un carácter común. Se representa con la letra N .
- MUESTRA: Conjunto de unidades de una población. Cuanto más significativa sea, mejor será la muestra. Se representa con la letra n .
- UNIDAD DE MUESTREO: Está formada por uno o más elementos de la población. El total de unidades de muestreo constituyen la población. Estas unidades son disjuntas entre sí y cada elemento de la población pertenece a una unidad de muestreo.
- PARÁMETRO: Es un resumen numérico de alguna variable observada de la población.
- ESTIMADOR: Un estimador de un parámetro, es un estadístico que se emplea para conocer el parámetro desconocido.
- ESTADÍSTICO: Es una función de los valores de la muestra. Es una variable aleatoria, cuyos valores dependen de la muestra seleccionada. Su distribución de probabilidad, se conoce como Distribución muestral del estadístico.
- ESTIMACIÓN: Este término indica que a partir de lo observado en una muestra (un resumen estadístico con las medidas que conocemos de Descriptiva) se extrapola o generaliza dicho resultado muestral a la población total, de modo que lo estimado es el valor generalizado a la población. Consiste en la búsqueda del valor de los parámetros poblacionales objeto de estudio. Puede ser puntual o por intervalo de confianza:

Estimación

- Puntual: cuando buscamos un valor concreto. Inferencia, estimación y contraste de hipótesis
- Intervalo de confianza: cuando determinamos un intervalo, dentro del cual se supone que va a estar el valor del parámetro que se busca con una cierta probabilidad.
- CONTRATE DE HIPÓTESIS: Consiste en determinar si es aceptable, partiendo de datos muestrales, que la característica o el parámetro poblacional estudiado tome un determinado valor o esté dentro de unos determinados valores.
- NIVEL DE CONFIANZA: Indica la proporción de veces que acertaríamos al afirmar que el parámetro está dentro del intervalo al seleccionar muchas muestras.

EL CONCEPTO DE ESTADÍSTICO Y DISTRIBUCIÓN MUESTRAL

- El objetivo de la inferencia es efectuar una generalización de los resultados de la muestra de la población. La tarea que nos ocupa ahora es conocer las distribuciones de la probabilidad de ciertas funciones de la muestra, es decir, variables aleatorias asociadas al muestreo o estadísticos muestrales. éstos serán útiles para hacer inferencia respecto a los parámetros desconocidos de una población.
- Por ello se habla de distribuciones muestrales, ya que están basados en el comportamiento de las muestras.
- El primer objetivo es conocer el concepto de distribución muestral de un estadístico; su comportamiento probabilístico dependerá del que tenga la variable X y del tamaño de las muestras.
- Sea una población donde se observa la variable aleatoria X . Esta variable X , tendrá una distribución de probabilidad, que puede ser conocida o desconocida, y ciertas características o parámetros poblacionales.
- El problema será encontrar una función que proporcione el mejor estimador de El estimador, T , del parámetro debe tener una distribución concentrada alrededor de la media y la varianza debe ser lo menor posible.

- Los estadísticos más usuales en inferencia y su distribución asociada considerando una población P sobre la que se estudia un carácter cuantitativo son:
- CONTRATE DE HIPÓTESIS: Consiste en determinar si es aceptable, partiendo de datos muestrales, que la característica o el parámetro poblacional estudiado tome un determinado valor o esté dentro de unos determinados valores.
- NIVEL DE CONFIANZA: Indica la proporción de veces que acertaríamos al afirmar que el parámetro está dentro del intervalo al seleccionar muchas muestras.

El objetivo de la inferencia es efectuar una generalización de los resultados de la muestra de la población. La tarea que nos ocupa ahora es conocer las distribuciones de la probabilidad de ciertas funciones de la muestra, es decir, variables aleatorias asociadas al muestreo o estadísticos muestrales. éstos serán útiles para hacer inferencia respecto a los parámetros desconocidos de una población. Por ello se habla de distribuciones muestrales, ya que están basados en el comportamiento de las muestras.

El primer objetivo es conocer el concepto de distribución muestral de un estadístico; su comportamiento probabilístico dependerá del que tenga la variable X y del tamaño de las muestras.

47.2 2. Pruebas de Hipótesis

47.2.1 2.1 Tipos de errores

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado
- el interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.
- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

Una prueba de hipótesis está formada por cinco partes

- La hipótesis nula, denotada por H_0 .
- La hipótesis alterativa, denorada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 33. Las dos hipótesis en competencias son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones:

- Rechazar H_0 y concluir que H_1 es verdadera.

- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 15. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 16. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 34. El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadística de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**.

Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**.

Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente.

Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 35. Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

47.3 2.2 Muestras grandes: una media poblacional

47.3.1 2.2.1 Cálculo de valor p

Definición 36. El **valor de p** (*p-value*) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.

Nota 41. Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.

Nota 42. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 37. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados

p	H_0	Significativa
$p < 0.01$	Rechazar	Altamente
$0.01 \leq p < 0.05$	Rechazar	Estadísticamente
$0.05 \leq p < 0.1$	No rechazar	Tendencia estadística
$0.1 \leq p$	No rechazar	No son estadísticamente

Nota 43. Para determinar el valor de p , encontrar el área en la cola después del estadístico de prueba. Si la prueba es de una cola, este es el valor de p . Si es de dos colas, éste valor encontrado es la mitad del valor de p . Rechazar H_0 cuando el valor de $p < \alpha$.

Hay dos tipos de errores al realizar una prueba de hipótesis

	H_0 es Verdadera	H_0 es Falsa
Rechazar H_0	Error tipo I	✓
Aceptar H_0	✓	Error tipo II

Definición 38. La probabilidad de cometer el error tipo II se define por β donde

$$\begin{aligned}\beta &= P\{\text{error tipo II}\} = P\{\text{Aceptar equivocadamente } H_0\} \\ &= P\{\text{Aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\}\end{aligned}$$

Nota 44. Cuando H_0 es falsa y H_1 es verdadera, no siempre es posible especificar un valor exacto de μ , sino más bien un rango de posibles valores. En lugar de arriesgarse a tomar una decisión incorrecta, es mejor concluir que no hay evidencia suficiente para rechazar H_0 , es decir en lugar de aceptar H_0 , no rechazar H_0 .

La bondad de una prueba estadística se mide por el tamaño de α y β , ambas deben de ser pequeñas. Una manera muy efectiva de medir la potencia de la prueba es calculando el complemento del error tipo II:

$$\begin{aligned}1 - \beta &= P\{\text{Rechazar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{\text{Rechazar } H_0 \text{ cuando } H_1 \text{ es verdadera}\}\end{aligned}$$

Definición 39. La **potencia de la prueba**, $1 - \beta$, mide la capacidad de que la prueba funciona como se necesita.

Ejemplo 17. La producción diaria de una planta química local ha promediado 880 toneladas en los últimos años. A la gerente de control de calidad le gustaría saber si este promedio ha cambiado en meses recientes. Ella selecciona al azar 50 días de la base de datos computarizada y calcula el promedio y la desviación estándar de las $n = 50$ producciones como $\bar{x} = 871$ toneladas y $s = 21$ toneladas, respectivamente. Pruebe la hipótesis apropiada usando $\alpha = 0.05$.

Solución 1. La hipótesis nula apropiada es:

$$\begin{aligned} H_0 &: \mu = 880 \\ &\text{y la hipótesis alternativa } H_1 \text{ es} \\ H_1 &: \mu \neq 880 \end{aligned}$$

el estimador puntual para μ es \bar{x} , entonces el estadístico de prueba es

$$\begin{aligned} z &= \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \\ &= \frac{871 - 880}{21/\sqrt{50}} = -3.03 \end{aligned}$$

Solución 2. Para esta prueba de dos colas, hay que determinar los dos valores de $z_{\alpha/2}$, es decir, $z_{\alpha/2} = \pm 1.96$, como $z > z_{\alpha/2}$, z cae en la zona de rechazo, por lo tanto la gerente puede rechazar la hipótesis nula y concluir que el promedio efectivamente ha cambiado. La probabilidad de rechazar H_0 cuando esta es verdadera es de 0.05.

Recordemos que el valor observado del estadístico de prueba es $z = -3.03$, la región de rechazo más pequeña que puede usarse y todavía seguir rechazando H_0 es $|z| > 3.03$, entonces $p = 2(0.012) = 0.0024$, que a su vez es menor que el nivel de significancia α asignado inicialmente, y además los resultados son **altamente significativos**.

Finalmente determinemos la potencia de la prueba cuando μ en realidad es igual a 870 toneladas.

Recordar que la región de aceptación está entre -1.96 y 1.96 , para $\mu = 880$, equivalentemente

$$874.18 < \bar{x} < 885.82$$

β es la probabilidad de aceptar H_0 cuando $\mu = 870$, calculemos los valores de z correspondientes a 874.18 y 885.82. Entonces

$$\begin{aligned} z_1 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{874.18 - 870}{21/\sqrt{50}} = 1.41 \\ z_1 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{885.82 - 870}{21/\sqrt{50}} = 5.33 \end{aligned}$$

por lo tanto

$$\begin{aligned} \beta &= P\{\text{aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{874.18 < \mu < 885.82 \text{ cuando } \mu = 870\} \\ &= P\{1.41 < z < 5.33\} = P\{1.41 < z\} \\ &= 1 - 0.9207 = 0.0793 \end{aligned}$$

entonces, la potencia de la prueba es

$$1 - \beta = 1 - 0.0793 = 0.9207$$

que es la probabilidad de rechazar correctamente H_0 cuando H_0 es falsa.

Determinar la potencia de la prueba para distintos valores de H_1 y graficarlos, *curva de potencia*

H_1	$(1 - \beta)$
865	
870	
872	
875	
877	
880	
883	
885	
888	
890	
895	

- Encontrar las regiones de rechazo para el estadístico z , para una prueba de
 - dos colas para $\alpha = 0.01, 0.05, 0.1$
 - una cola superior para $\alpha = 0.01, 0.05, 0.1$
 - una cola inferior para $\alpha = 0.01, 0.05, 0.1$
- Suponga que el valor del estadístico de prueba es
 - $z = -2.41$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = 2.16$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = 1.15$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = -2.78$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = -1.81$, sacar las conclusiones correspondientes para los incisos anteriores.
- Encuentre el valor de p para las pruebas de hipótesis correspondientes a los valores de z del ejercicio anterior.
- Para las pruebas dadas en el ejercicio 2, utilice el valor de p , determinado en el ejercicio 3, para determinar la significancia de los resultados.
- Una muestra aleatoria de $n = 45$ observaciones de una población con media $\bar{x} = 2.4$, y desviación estándar $s = 0.29$. Suponga que el objetivo es demostrar que la media poblacional μ excede 2.3.
 - Defina la hipótesis nula y alternativa para la prueba.
 - Determine la región de rechazo para un nivel de significancia de: $\alpha = 0.1, 0.05, 0.01$.
 - Determine el error estándar de la media muestral.
 - Calcule el valor de p para los estadísticos de prueba definidos en los incisos anteriores.
 - Utilice el valor de p para sacar una conclusión al nivel de significancia α .
 - Determine el valor de β cuando $\mu = 2.5$
 - Graficar la curva de potencia para la prueba.

47.3.2 2.2.2 Prueba de hipótesis para la diferencia entre dos medias poblacionales

El estadístico que resume la información muestral respecto a la diferencia en medias poblacionales $(\mu_1 - \mu_2)$ es la diferencia de las medias muestrales $(\bar{x}_1 - \bar{x}_2)$, por tanto al probar la diferencia entre las medias muestrales se verifica que la diferencia real entre las medias poblacionales difiere de un valor especificado, $(\mu_1 - \mu_2) = D_0$, se puede usar el error estándar de $(\bar{x}_1 - \bar{x}_2)$, es decir

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

cuyo estimador está dado por

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

El procedimiento para muestras grandes es:

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

- 2) **Hipótesis Alternativa**

Prueba de una Cola	Prueba de dos colas
$H_1 : (\mu_1 - \mu_2) > D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$
$H_1 : (\mu_1 - \mu_2) < D_0$	

- 3) Estadístico de prueba:

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_0$	
$z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cuando $p < \alpha$	

Ejemplo 18. Para determinar si ser propietario de un automóvil afecta el rendimiento académico de un estudiante, se tomaron dos muestras aleatorias de 100 estudiantes varones. El promedio de calificaciones para los $n_1 = 100$ no propietarios de un auto tuvieron un promedio y varianza de $\bar{x}_1 = 2.7$ y $s_1^2 = 0.36$, respectivamente, mientras que para la segunda muestra con $n_2 = 100$ propietarios de un auto, se tiene $\bar{x}_2 = 2.54$ y $s_2^2 = 0.4$. Los datos presentan suficiente evidencia para indicar una diferencia en la media en el rendimiento académico entre propietarios y no propietarios de un automóvil? Hacer pruebas para $\alpha = 0.01, 0.05$ y $\alpha = 0.1$.

Solución 3. – Solución utilizando la técnica de regiones de rechazo: realizando las operaciones $z = 1.84$, determinar si excede los valores de $z_{\alpha/2}$.

- Solución utilizando el p -value: Calcular el valor de p , la probabilidad de que z sea mayor que $z = 1.84$ o menor que $z = -1.84$, se tiene que $p = 0.0658$. Concluir.
- Si el intervalo de confianza que se construye contiene el valor del parámetro especificado por H_0 , entonces ese valor es uno de los posibles valores del parámetro y H_0 no debe ser rechazada.
- Si el valor hipotético se encuentra fuera de los límites de confianza, la hipótesis nula es rechazada al nivel de significancia α .

1. Del libro Mendenhall resolver los ejercicios 9.18, 9.19 y 9.20(Mendenhall).
2. Del libro Mendenhall resolver los ejercicios: 9.23, 9.26 y 9.28.

47.3.3 2.2.3 Prueba de Hipótesis para una Proporción Binomial

Para una muestra aleatoria de n intentos idénticos, de una población binomial, la proporción muestral \hat{p} tiene una distribución aproximadamente normal cuando n es grande, con media p y error estándar

$$SE = \sqrt{\frac{pq}{n}}.$$

La prueba de hipótesis de la forma

$$\begin{aligned} H_0 &: p = p_0 \\ H_1 &: p > p_0, \text{ o } p < p_0 \text{ o } p \neq p_0 \end{aligned}$$

El estadístico de prueba se construye con el mejor estimador de la proporción verdadera, \hat{p} , con el estadístico de prueba z , que se distribuye normal estándar. El procedimiento es

- 1) Hipótesis nula: $H_0 : p = p_0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis alternativa	$H_1 : p > p_0$ $H_1 : p < p_0$	$p \neq p_0$

- 3) Estadístico de prueba:

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{pq}{n}}}, \hat{p} = \frac{x}{n}$$

donde x es el número de éxitos en n intentos binomiales.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_0$ $z < -z_\alpha$ cuando $H_1 : p < p_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Ejemplo 19. A cualquier edad, alrededor del 20% de los adultos de cierto país realiza actividades de acondicionamiento físico al menos dos veces por semana. En una encuesta local de $n = 100$ adultos de más de 40 a nos, un total de 15 personas indicaron que realizaron actividad física al menos dos veces por semana. Estos datos indican que el porcentaje de participación para adultos de más de 40 a nos de edad es considerablemente menor a la cifra del 20%? Calcule el valor de p y úselo para sacar las conclusiones apropiadas.

1. Resolver los ejercicios: 9.30, 9.32, 9.33, 9.35 y 9.39.

47.3.4 2.2.4 Prueba de Hipótesis diferencia entre dos Proporciones Binomiales

Nota 45. Cuando se tienen dos muestras aleatorias independientes de dos poblaciones binomiales, el objetivo del experimento puede ser la diferencia $(p_1 - p_2)$ en las proporciones de individuos u objetos que poseen una característica específica en las dos poblaciones. En este caso se pueden utilizar los estimadores de las dos proporciones $(\hat{p}_1 - \hat{p}_2)$ con error estándar dado por

$$SE = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

considerando el estadístico z con un nivel de significancia $(1 - \alpha) 100\%$

Nota 46. La hipótesis nula a probarse es de la forma

$H_0: p_1 = p_2$ o equivalentemente $(p_1 - p_2) = 0$, contra una hipótesis alternativa H_1 de una o dos colas.

Nota 47. Para estimar el error estándar del estadístico z , se debe de utilizar el hecho de que suponiendo que H_0 es verdadera, las dos proporciones son iguales a algún valor común, p . Para obtener el mejor estimador de p es

$$p = \frac{\text{número total de éxitos}}{\text{Número total de pruebas}} = \frac{x_1 + x_2}{n_1 + n_2}$$

1) **Hipótesis Nula:** $H_0 : (p_1 - p_2) = 0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 :$	$H_1 : (p_1 - p_2) > 0$ $H_1 : (p_1 - p_2) < 0$	$H_1 : (p_1 - p_2) \neq 0$

3) Estadístico de prueba:

$$z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{pq}{n_1} + \frac{pq}{n_2}}}$$

donde $\hat{p}_1 = x_1/n_1$ y $\hat{p}_2 = x_2/n_2$, dado que el valor común para p_1 y p_2 es p , entonces $\hat{p} = \frac{x_1+x_2}{n_1+n_2}$ y por tanto el estadístico de prueba es

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_\alpha$ $z < -z_\alpha$ cuando $H_1 : p < p_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Ejemplo 20. Los registros de un hospital, indican que 52 hombres de una muestra de 1000 contra 23 mujeres de una muestra de 1000 fueron ingresados por enfermedad del corazón. Estos datos presentan suficiente evidencia para indicar un porcentaje más alto de enfermedades del corazón entre hombres ingresados al hospital?, utilizar distintos niveles de confianza de α .

1. Resolver los ejercicios 9.42
2. Resolver los ejercicios: 9.45, 9.48, 9.50

47.4 2.3 Muestras Pequeñas

47.4.1 2.3.1 Una media poblacional

1) **Hipótesis Nula:** $H_0 : \mu = \mu_0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 :$	$H_1 : \mu > \mu_0$ $H_1 : \mu < \mu_0$	$H_1 : \mu \neq \mu_0$

3) Estadístico de prueba:

$$t = \frac{\bar{x} - \mu_0}{\sqrt{\frac{s^2}{n}}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Ejemplo 21. Las etiquetas en latas de un galón de pintura por lo general indican el tiempo de secado y el área puede cubrir una capa. Casi todas las marcas de pintura indican que, en una capa, un galón cubrirá entre 250 y 500 pies cuadrados, dependiendo de la textura de la superficie a pintarse, un fabricante, sin embargo afirma que un galón de su pintura cubrirá 400 pies cuadrados de área superficial. Para probar su afirmación, una muestra aleatoria de 10 latas de un galón de pintura blanca se empleó para pintar 10 áreas idénticas usando la misma clase de equipo. Las áreas reales en pies cuadrados cubiertas por estos 10 galones de pintura se dan a continuación:

310	311	412	368	447
376	303	410	365	350

Ejemplo 22. Los datos presentan suficiente evidencia para indicar que el promedio de la cobertura difiere de 400 pies cuadrados? encuentre el valor de p para la prueba y úselo para evaluar la significancia de los resultados.

1. Resolver los ejercicios: 10.2, 10.3, 10.5, 10.7, 10.9, 10.13 y 10.16

47.4.2 2.3.2 Diferencia entre dos medias poblacionales: M.A.I.

Nota 48. Cuando los tamaños de muestra son pequeños, no se puede asegurar que las medias muestrales sean normales, pero si las poblaciones originales son normales, entonces la distribución muestral de la diferencia de las medias muestrales, $(\bar{x}_1 - \bar{x}_2)$, será normal con media $(\mu_1 - \mu_2)$ y error estándar

$$ES = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\mu_1 - \mu_2) > D_0$ $H_1 : (\mu_1 - \mu_2) < D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$

- 3) Estadístico de prueba:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

donde

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

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_0$ $z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

47.4.3 2.3.3 Diferencia entre dos medias poblacionales: Diferencias Pareadas

- 1) **Hipótesis Nula:** $H_0 : \mu_d = 0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 : \mu_d$	$H_1 : \mu_d > 0$ $H_1 : \mu_d < 0$	$H_1 : \mu_d \neq 0$

- 3) Estadístico de prueba:

$$t = \frac{\bar{d}}{\sqrt{\frac{s_d^2}{n}}}$$

donde n es el número de diferencias pareadas, \bar{d} es la media de las diferencias muestrales, y s_d es la desviación estándar de las diferencias muestrales.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

47.4.4 2.3.4 Inferencias con respecto a la Varianza Poblacional

- 1) **Hipótesis Nula:** $H_0 : \sigma^2 = \sigma_0^2$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: H_1	$H_1 : \sigma^2 > \sigma_0^2$ $H_1 : \sigma^2 < \sigma_0^2$	$H_1 : \sigma^2 \neq \sigma_0^2$

- 3) Estadístico de prueba:

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$\chi^2 > \chi_\alpha^2$ $\chi^2 < \chi_{(1-\alpha)}^2$ cuando $H_1 : \chi^2 < \chi_0^2$ cuando $p < \alpha$	$\chi^2 > \chi_{\alpha/2}^2$ o $\chi^2 < \chi_{(1-\alpha/2)}^2$

Los valores críticos de χ^2 , están basados en $(n_1 +)$ grados de libertad.

47.4.5 2.3.5 Comparación de dos varianzas poblacionales

- 1) **Hipótesis Nula** $H_0 : (\sigma_1^2 - \sigma_2^2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\sigma_1^2 - \sigma_2^2) > D_0$ $H_1 : (\sigma_1^2 - \sigma_2^2) < D_0$	$H_1 : (\sigma_1^2 - \sigma_2^2) \neq D_0$

- 3) Estadístico de prueba:

$$F = \frac{s_1^2}{s_2^2}$$

donde s_1^2 es la varianza muestral más grande.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$F > F_\alpha$ cuando $p < \alpha$	$F > F_{\alpha/2}$

47.5 3. Análisis de Regresion Lineal (RL)

Nota 49. — En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.

47.5.1 3.1 Regresión Lineal Simple (RLS)

- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (47.1)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

47.5.2 3.2 Método de Mínimos Cuadrados

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (47.2)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación ?? se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (47.3)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (47.4)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (47.5)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (47.6)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (47.7)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (47.8)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (47.9)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (47.10)$$

47.5.3 3.3 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Nota 50. – Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo.

– Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right) \\ &= \frac{1}{S_{xx}} E\left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x})\right) \\ &= \frac{1}{S_{xx}} \left[\beta_0 E\left(\sum_{k=1}^n (x_k - \bar{x})\right) + E\left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x})\right) \right. \\ &\quad \left. + E\left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x})\right) \right] = \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1 \end{aligned}$$

por lo tanto

$$E(\hat{\beta}_1) = \beta_1 \quad (47.11)$$

Nota 51. Es decir, $\hat{\beta}_1$ es un estimador insesgado.

Ahora calculemos la varianza:

$$\begin{aligned} V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) \\ &= \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) = \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 \\ &= \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}} \end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad (47.12)$$

Proposición 11.

$$\begin{aligned} E(\hat{\beta}_0) &= \beta_0, \\ V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \\ Cov(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{\sigma^2 \bar{x}}{S_{xx}}. \end{aligned}$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$$e_k = y_k - \hat{y}_k, \text{ se le denomina } \mathbf{residuo}.$$

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \quad (47.13)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$\begin{aligned} SC_E &= \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \\ E(SC_E) &= (n-2)\sigma^2, \text{ por lo tanto} \\ \hat{\sigma}^2 &= \frac{SC_E}{n-2} = MC_E \text{ es un estimador insesgado de } \sigma^2. \end{aligned}$$

47.5.4 3.4 Prueba de Hipótesis en RLS

- Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza.
- Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$.

Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_1, \sigma^2)$. De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??).

Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (47.14)$$

que se distribuye t con $n-2$ grados de libertad bajo $H_0: \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (47.15)$$

Para β_0 se puede proceder de manera análoga para

$$H_0: \beta_0 = \beta_{0,0},$$

$$H_1: \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}, \quad (47.16)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (47.17)$$

- No rechazar $H_0 : \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0 : \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0 : \beta_1 = 0$ puede realizarse de la siguiente manera:

$$\begin{aligned} S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 \\ &= \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\ &= \sum_{k=1}^n \left[(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2 \right] \\ &= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ &\quad \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) \\ &= \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\ &= \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &= \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &= \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0. \end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (47.18)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (47.19)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (47.20)$$

Por lo tanto la ecuación (??) se puede reescribir como

$$S_{yy} = SC_R + SC_E \quad (47.21)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$

$$\begin{aligned} S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}) \\ S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R \end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 12.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx} \quad (47.22)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E}$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$.

El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n - 2$	MC_E	
Total	S_{yy}	$n - 1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}} \quad (47.23)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

47.5.5 Estimación de Intervalos en RLS

- Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros.
- El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión.

Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n - 2$ grados de libertad. Por tanto un intervalo de confianza de $100(1 - \alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}}. \quad (47.24)$$

De igual manera, para β_0 un intervalo de confianza al $100(1 - \alpha)\%$ es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (47.25)$$

47.5.6 Predicción

Supongamos que se tiene un valor x_0 de interés, entonces la estimación puntual de este nuevo valor

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 \quad (47.26)$$

Nota 52. Esta nueva observación es independiente de las utilizadas para obtener el modelo de regresión, por tanto, el intervalo en torno a la recta de regresión es inapropiado, puesto que se basa únicamente en los datos empleados para ajustar el modelo de regresión.

El intervalo de confianza en torno a la recta de regresión se refiere a la respuesta media verdadera $x = x_0$, no a observaciones futuras.

Sea y_0 la observación futura en $x = x_0$, y sea \hat{y}_0 dada en la ecuación anterior, el estimador de y_0 . Si se define la variable aleatoria

$$w = y_0 - \hat{y}_0,$$

esta se distribuye normalmente con media cero y varianza

$$V(w) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]$$

dado que y_0 es independiente de \hat{y}_0 , por lo tanto el intervalo de predicción al nivel α para futuras observaciones x_0 es

$$\begin{aligned} \hat{y}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]} &\leq y_0 \\ &\leq \hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]}. \end{aligned}$$

47.5.7 Coeficiente de Determinación

La cantidad

$$R^2 = \frac{SC_R}{S_{yy}} = 1 - \frac{SC_E}{S_{yy}} \quad (47.27)$$

se denomina coeficiente de determinación y se utiliza para saber si el modelo de regresión es suficiente o no. Se puede demostrar que $0 \leq R^2 \leq 1$, una manera de interpretar este valor es que si $R^2 = k$, entonces el modelo de regresión explica el $k * 100\%$ de la variabilidad en los datos.

R^2

- No mide la magnitud de la pendiente de la recta de regresión
- Un valor grande de R^2 no implica una pendiente empinada.
- No mide la suficiencia del modelo.
- Valores grandes de R^2 no implican necesariamente que el modelo de regresión proporcionará predicciones precisas para futuras observaciones.

47.6 Análisis de Varianza

Para analizar el ajuste de regresión se utiliza el método de **Análisis de Varianza (ANOVA)**, el en cuál se estudia la variación de la variable dependiente, subdividiéndola en dos componentes significativos. Recordemos las ecuaciones ?? y ??:

$$S_{yy} = SC_R + SC_E.$$

SC_R Se le denomina **suma de cuadrados de la regresión** y refleja la cantidad de variación de los valores de y que es explicada por el modelo, para nuestro caso: la recta propuesta.

SC_E Se le denomina suma de cuadrados del error, que es la variación o diferencia que hay entre los valores originales y los obtenidos mediante el ajuste.

De lo anterior se desprende que estamos interesados en validar nuestro modelo dado en la ecuación (??), es decir,

$$y = \beta_0 + \beta_1 x + \epsilon$$

que en realidad el parámetro β_1 ha sido bien estimado:

Supongamos que se desea probar la hipótesis

$$H_0 : \beta_1 = 0$$

$$H_1 : \beta_1 \neq 0$$

Donde la hipótesis nula nos dice que el modelo en realidad debe de ser: $\mu_{Y|x} = \beta_0$, es decir, las variaciones en los valores de Y son independientes de los valores de x . Se puede demostrar que bajo la hipótesis nula los términos

- SC_R/σ^2 se distribuye χ^2 con 1 grado de libertad
- SC_E/σ^2 se distribuye χ^2 con $n - 2$ grado de libertad.

e independientes, y por tanto S_{yy} , también llamada **suma total de cuadrados corregida: STCC**, se distribuye χ^2 con $n - 1$ grados de libertad.

Para realizar esta prueba de hipótesis se calcula el cociente

$$f = \frac{SC_R/1}{SC_E/(n-2)} = \frac{SC_R}{s^2}$$

y se rechaza H_0 a un nivel de significancia α si $f > f_\alpha(1, (n-2))$, esto se puede realizar mediante una tabla, llamada tabla de análisis de varianza, cuando a las distintas sumas de cuadrados se les divide por sus grados de libertad, se les denomina **cuadrados medios**.

47.7 3. Análisis de Regresión Lineal (RL)

Nota 53. — En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.

47.7.1 3.1 Regresión Lineal Simple (RLS)

- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (47.28)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

47.7.2 3.2 Método de Mínimos Cuadrados

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (47.29)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación ?? se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (47.30)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (47.31)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (47.32)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (47.33)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (47.34)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (47.35)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (47.36)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (47.37)$$

47.7.3 3.3 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Nota 54. – Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo.

- Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$\begin{aligned}
 E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right) \\
 &= \frac{1}{S_{xx}} E\left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x})\right) \\
 &= \frac{1}{S_{xx}} \left[\beta_0 E\left(\sum_{k=1}^n (x_k - \bar{x})\right) + E\left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x})\right) \right. \\
 &\quad \left. + E\left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x})\right) \right] = \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1
 \end{aligned}$$

por lo tanto

$$E(\hat{\beta}_1) = \beta_1 \quad (47.38)$$

Nota 55. Es decir, $\hat{\beta}_1$ es un estimador insesgado.

Ahora calculemos la varianza:

$$\begin{aligned}
 V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) \\
 &= \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) = \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 \\
 &= \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}}
 \end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad (47.39)$$

Proposición 13.

$$\begin{aligned}
 E(\hat{\beta}_0) &= \beta_0, \\
 V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \\
 Cov(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{\sigma^2 \bar{x}}{S_{xx}}.
 \end{aligned}$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$$e_k = y_k - \hat{y}_k, \text{ se le denomina } \mathbf{residuo}.$$

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \quad (47.40)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$\begin{aligned} SC_E &= \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \\ E(SC_E) &= (n-2)\sigma^2, \text{ por lo tanto} \\ \hat{\sigma}^2 &= \frac{SC_E}{n-2} = MC_E \text{ es un estimador insesgado de } \sigma^2. \end{aligned}$$

47.7.4 3.4 Prueba de Hipótesis en RLS

- Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza.
- Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$.

Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_1, \sigma^2)$. De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??).

Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (47.41)$$

que se distribuye t con $n-2$ grados de libertad bajo $H_0: \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (47.42)$$

Para β_0 se puede proceder de manera análoga para

$$H_0: \beta_0 = \beta_{0,0},$$

$$H_1: \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}, \quad (47.43)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (47.44)$$

- No rechazar $H_0 : \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0 : \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0 : \beta_1 = 0$ puede realizarse de la siguiente manera:

$$\begin{aligned} S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 \\ &= \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\ &= \sum_{k=1}^n \left[(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2 \right] \\ &= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ &\quad \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) \\ &= \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\ &= \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &= \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &= \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0. \end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (47.45)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (47.46)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (47.47)$$

Por lo tanto la ecuación (??) se puede reescribir como

$$S_{yy} = SC_R + SC_E \quad (47.48)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$

$$\begin{aligned} S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}) \\ S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R \end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 14.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx} \quad (47.49)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E}$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$.

El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n - 2$	MC_E	
Total	S_{yy}	$n - 1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}} \quad (47.50)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

47.7.5 Estimación de Intervalos en RLS

- Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros.
- El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión.

Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n - 2$ grados de libertad. Por tanto un intervalo de confianza de $100(1 - \alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}}. \quad (47.51)$$

De igual manera, para β_0 un intervalo de confianza al $100(1 - \alpha)\%$ es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (47.52)$$

47.7.6 Predicción

Supongamos que se tiene un valor x_0 de interés, entonces la estimación puntual de este nuevo valor

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 \quad (47.53)$$

Nota 56. Esta nueva observación es independiente de las utilizadas para obtener el modelo de regresión, por tanto, el intervalo en torno a la recta de regresión es inapropiado, puesto que se basa únicamente en los datos empleados para ajustar el modelo de regresión.

El intervalo de confianza en torno a la recta de regresión se refiere a la respuesta media verdadera $x = x_0$, no a observaciones futuras.

Sea y_0 la observación futura en $x = x_0$, y sea \hat{y}_0 dada en la ecuación anterior, el estimador de y_0 . Si se define la variable aleatoria

$$w = y_0 - \hat{y}_0,$$

esta se distribuye normalmente con media cero y varianza

$$V(w) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]$$

dado que y_0 es independiente de \hat{y}_0 , por lo tanto el intervalo de predicción al nivel α para futuras observaciones x_0 es

$$\begin{aligned} \hat{y}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]} &\leq y_0 \\ &\leq \hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]}. \end{aligned}$$

47.7.7 Coeficiente de Determinación

La cantidad

$$R^2 = \frac{SC_R}{S_{yy}} = 1 - \frac{SC_E}{S_{yy}} \quad (47.54)$$

se denomina coeficiente de determinación y se utiliza para saber si el modelo de regresión es suficiente o no. Se puede demostrar que $0 \leq R^2 \leq 1$, una manera de interpretar este valor es que si $R^2 = k$, entonces el modelo de regresión explica el $k * 100\%$ de la variabilidad en los datos.

R^2

- No mide la magnitud de la pendiente de la recta de regresión
- Un valor grande de R^2 no implica una pendiente empinada.
- No mide la suficiencia del modelo.
- Valores grandes de R^2 no implican necesariamente que el modelo de regresión proporcionará predicciones precisas para futuras observaciones.

47.8 Análisis de Varianza

Para analizar el ajuste de regresión se utiliza el método de **Análisis de Varianza (ANOVA)**, el en cuál se estudia la variación de la variable dependiente, subdividiéndola en dos componentes significativos. Recordemos las ecuaciones ?? y ??:

$$S_{yy} = SC_R + SC_E.$$

SC_R Se le denomina **suma de cuadrados de la regresión** y refleja la cantidad de variación de los valores de y que es explicada por el modelo, para nuestro caso: la recta propuesta.

SC_E Se le denomina suma de cuadrados del error, que es la variación o diferencia que hay entre los valores originales y los obtenidos mediante el ajuste.

De lo anterior se desprende que estamos interesados en validar nuestro modelo dado en la ecuación (??), es decir,

$$y = \beta_0 + \beta_1 x + \epsilon$$

que en realidad el parámetro β_1 ha sido bien estimado:

Supongamos que se desea probar la hipótesis

$$H_0 : \beta_1 = 0$$

$$H_1 : \beta_1 \neq 0$$

Donde la hipótesis nula nos dice que el modelo en realidad debe de ser: $\mu_{Y|x} = \beta_0$, es decir, las variaciones en los valores de Y son independientes de los valores de x . Se puede demostrar que bajo la hipótesis nula los términos

- SC_R/σ^2 se distribuye χ^2 con 1 grado de libertad
- SC_E/σ^2 se distribuye χ^2 con $n - 2$ grado de libertad.

Fuente de Variación	Suma de Cuadrados	Grados de Libertad	Cuadrado Medio	f calculada
Regresión	SC_R	1	$\frac{SC_R}{1}$	$\frac{SC_R}{s^2}$
Error	SC_E	$n - 2$	$s^2 = \frac{SC_E}{n-2}$	
Total	$STCC$	$n - 1$		

Table 47.1: Análisis de Varianza para la prueba $\beta_1 = 0$

e independientes, y por tanto S_{yy} , también llamada **suma total de cuadrados corregida: STCC**, se distribuye χ^2 con $n - 1$ grados de libertad.

Para realizar esta prueba de hipótesis se calcula el cociente

$$f = \frac{SC_R/1}{SC_E/(n-2)} = \frac{SC_R}{s^2}$$

y se rechaza H_0 a un nivel de significancia α si $f > f_\alpha(1, (n-2))$, esto se puede realizar mediante una tabla, llamada tabla de análisis de varianza, cuando a las distintas sumas de cuadrados se les divide por sus grados de libertad, se les denomina **cuadrados medios**.

Se rechaza la hipótesis nula, cuando el estadístico F calculado excede al valor crítico $f_\alpha(1 - n - 2)$, y entonces se concluye que existe evidencia sobre la variación respecto al modelo ajustado. Si el estadístico F está en la región de no rechazo, se concluye que los datos no reflejan evidencia suficiente para sostener que el modelo ajustado.

Para hacer la prueba de hipótesis

$$H_0 : \beta_1 = \beta_{10}$$

$$H_1 : \beta_1 \neq \beta_{10}$$

se utiliza el estadístico:

$$T = \frac{B_1 - \beta_{10}}{S/\sqrt{S_{xx}}}$$

donde T se distribuye t con $n - 2$ grados de libertad. La hipótesis se rechaza si $|t| > t_{\alpha/2}$ con un nivel de confianza α .

Nota 57. Para el caso en que $\beta_{10} = 0$, se tiene que el valor del estadístico se convierte en

$$T = \frac{b_1 - \beta_{10}}{s/\sqrt{S_{xx}}}$$

y entonces el análisis es similar al dado en la tabla ??, y lo que se está diciendo es que la variación depende totalmente del azar.

Nota 58. El Análisis de Varianza utiliza la distribución F en lugar de la distribución t .

Supongamos que se tienen observaciones repetidas de las respuestas para k valores distintos de x , es decir: para x_1, x_2, \dots, x_k se tienen $y_{1,1}, y_{1,2}, \dots, y_{1,n_1}$ valores observados para la variable aleatoria Y_1 , $y_{2,1}, y_{2,2}, \dots, y_{2,n_2}$ valores observados para la variable aleatoria Y_2 , y así sucesivamente para $y_{k,1}, y_{k,2}, \dots, y_{k,n_k}$ valores observados para la variable aleatoria Y_k , de tal manera que

$$n = \sum_{i=1}^k n_i$$

$$Y = \begin{bmatrix} y_{1,1} & y_{2,1} & \cdots & Y_{k,1} \\ y_{1,2} & y_{2,2} & \cdots & Y_{k,2} \\ \vdots & \vdots & \vdots & \vdots \\ y_{1,j} & y_{2,j} & y_{i,j} & y_{k,j} \\ \vdots & \vdots & \vdots & \vdots \\ y_{1,n_1} & y_{2,n_2} & \cdots & Y_{k,n_k} \end{bmatrix}$$

entonces, si definimos $y_i = T_i = \sum_{j=1}^{n_i} y_{i,j}$, se tiene que $\bar{y}_i = \frac{T_i}{n_i}$ Cómo se ve la matriz para el caso en que:

- $n_4 = 3$ mediciones de Y
- Simular en R, para los casos en que $n_1 = 4$, $n_2 = 6$, $n_3 = 5$, y $n_4 = 8$,

La suma de cuadrados del error se divide en dos partes: la cantidad debida a la variación entre los valores de Y para los valores dados de x , y lo que se denomina **falta de ajuste** que es una medida de la variación sistemática introducida por los términos de orden superior. Para nuestro caso en específico, estos son términos de x distintos de la contribución lineal de primer orden.

Hasta el momento, dado que hemos considerado un modelo lineal, se asume que este segundo componente no existe, y por tanto la suma de cuadrados de error depende totalmente de los errores aleatorios. En consecuencia tenemos que $s^2 = \frac{SCE}{(n-2)}$ es un estimador insesgado para σ^2 . Sin embargo, si el modelo no ajusta correctamente a los datos, lo que tenemos es una sobre estimación del valor de σ^2 y por tanto será un estimador sesgado del mismo.

Para obtener un estimador insesgado se calcula

$$s^2 = \frac{\sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}{n_i - 1}, \text{ para } i = 1, 2, \dots, k$$

después de hacer unas operaciones se puede obtener:

$$s^2 = \frac{\sum_{i=1}^k (n_i - 1) s_i^2}{n_i - k} = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2}{n - k}.$$

El numerador de s^2 es una medida del **error experimental puro** o **falta de ajuste**

Para determinar el cuadrado del error en: error puro y la falta de ajuste:

- Se calcula la suma de cuadrados del error puro:

$$\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$$

esta suma de cuadrados tiene $n - k$ grados de libertad, y el cuadrado medio resultante es el estimador insesgado s^2 de σ^2 .

- Restar la suma de cuadrados del error puro de la suma de cuadrados del error, SCE, resultando la suma de cuadrados por ajuste. Los grados de libertad de la falta de ajuste se obtienen por la resta:
 $(n - 2) - (n - k) = k - 2$.

La prueba de hipótesis en un problema de regresión con mediciones repetidas de la respuesta se ilustra en la tabla ??:

Fuente de Variación	Suma de Cuadrados	Grados de Libertad	Cuadrado Medio	f calculada
Regresión	SC_R	1	SC_R	$\frac{SC_R}{s^2}$
Error	SC_E	$n - 2$		
Falta de ajuste	$SCE - SCE(puro)$	$k - 2$	$\frac{SCE - SCE(puro)}{k - 2}$	$\frac{SCE - SCE(puro)}{s^2(k - 2)}$
Error Puro	$SCE(puro)$	$n - k$	$s^2 = \frac{SCE(puro)}{n - k}$	
Total	$STCC$	$n - 1$		

Table 47.2: Análisis de Varianza para la prueba $\beta_1 = 0$

47.9 Pruebas de Hipótesis

47.9.1 Tipos de errores

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado
- el interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.
- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

Una prueba de hipótesis está formada por cinco partes

- La hipótesis nula, denotada por H_0 .
- La hipótesis alterativa, denotada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 40. Las dos hipótesis en competencias son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones

- Rechazar H_0 y concluir que H_1 es verdadera.
- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 23. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 24. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 41. *El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadística de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.*

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**.

Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**.

Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente.

Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 42. *Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es*

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

47.10 Muestras grandes: una media poblacional

47.10.1 Cálculo de valor p

Definición 43. *El **valor de p** (p -value) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.*

Nota 59. *Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.*

Nota 60. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 44. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados

p	H_0	Significativa
$p < 0.01$	Rechazar	Altamente
$0.01 \leq p < 0.05$	Rechazar	Estadísticamente
$0.05 \leq p < 0.1$	No rechazar	Tendencia estadística
$0.1 \leq p$	No rechazar	No son estadísticamente

Nota 61. Para determinar el valor de p , encontrar el área en la cola después del estadístico de prueba. Si la prueba es de una cola, este es el valor de p . Si es de dos colas, éste valor encontrado es la mitad del valor de p . Rechazar H_0 cuando el valor de $p < \alpha$.

Hay dos tipos de errores al realizar una prueba de hipótesis

	H_0 es Verdadera	H_0 es Falsa
Rechazar H_0	Error tipo I	✓
Aceptar H_0	✓	Error tipo II

Definición 45. La probabilidad de cometer el error tipo II se define por β donde

$$\begin{aligned}\beta &= P\{\text{error tipo II}\} = P\{\text{Aceptar equivocadamente } H_0\} \\ &= P\{\text{Aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\}\end{aligned}$$

Nota 62. Cuando H_0 es falsa y H_1 es verdadera, no siempre es posible especificar un valor exacto de μ , sino más bien un rango de posibles valores. En lugar de arriesgarse a tomar una decisión incorrecta, es mejor concluir que no hay evidencia suficiente para rechazar H_0 , es decir en lugar de aceptar H_0 , no rechazar H_0 .

La bondad de una prueba estadística se mide por el tamaño de α y β , ambas deben de ser pequeñas. Una manera muy efectiva de medir la potencia de la prueba es calculando el complemento del error tipo II:

$$\begin{aligned}1 - \beta &= P\{\text{Rechazar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{\text{Rechazar } H_0 \text{ cuando } H_1 \text{ es verdadera}\}\end{aligned}$$

Definición 46. La **potencia de la prueba**, $1 - \beta$, mide la capacidad de que la prueba funciona como se necesita.

Ejemplo 25. La producción diaria de una planta química local ha promediado 880 toneladas en los últimos años. A la gerente de control de calidad le gustaría saber si este promedio ha cambiado en meses recientes. Ella selecciona al azar 50 días de la base de datos computarizada y calcula el promedio y la desviación estándar de las $n = 50$ producciones como $\bar{x} = 871$ toneladas y $s = 21$ toneladas, respectivamente. Pruebe la hipótesis apropiada usando $\alpha = 0.05$.

Solución 4. La hipótesis nula apropiada es:

$$H_0 : \mu = 880$$

y la hipótesis alternativa H_1 es

$$H_1 : \mu \neq 880$$

el estimador puntual para μ es \bar{x} , entonces el estadístico de prueba es

$$\begin{aligned} z &= \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \\ &= \frac{871 - 880}{21/\sqrt{50}} = -3.03 \end{aligned}$$

Solución 5. Para esta prueba de dos colas, hay que determinar los dos valores de $z_{\alpha/2}$, es decir, $z_{\alpha/2} = \pm 1.96$

47.11 Pruebas de Hipótesis

47.11.1 Tipos de errores

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado
- el interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.
- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.
- La decisión de aceptar o rechazar la hipótesis nula se basa en un estadístico calculado a partir de la muestra. Esto necesariamente implica la existencia de un error.

47.12 Pruebas de Hipótesis

47.12.1 Tipos de errores

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado
- el interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.

- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

Una prueba de hipótesis está formada por cinco partes

- La hipótesis nula, denotada por H_0 .
- La hipótesis alterativa, denorada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 47. Las dos hipótesis en competencias son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones

- Rechazar H_0 y concluir que H_1 es verdadera.
- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 26. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 27. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 48. El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadística de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**.

Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**.

Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente.

Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 49. Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

47.13 Muestras grandes: una media poblacional

47.13.1 Cálculo de valor p

Definición 50. El **valor de p** (*p-value*) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.

Nota 63. Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.

Nota 64. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 51. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados

p	H_0	Significativa
$p < 0.01$	Rechazar	Altamente
$0.01 \leq p < 0.05$	Rechazar	Estadísticamente
$0.05 \leq p < 0.1$	No rechazar	Tendencia estadística
$0.1 \leq p$	No rechazar	No son estadísticamente

Nota 65. Para determinar el valor de p , encontrar el área en la cola después del estadístico de prueba. Si la prueba es de una cola, este es el valor de p . Si es de dos colas, éste valor encontrado es la mitad del valor de p . Rechazar H_0 cuando el valor de $p < \alpha$.

Hay dos tipos de errores al realizar una prueba de hipótesis

	H_0 es Verdadera	H_0 es Falsa
Rechazar H_0	Error tipo I	✓
Aceptar H_0	✓	Error tipo II

Definición 52. La probabilidad de cometer el error tipo II se define por β donde

$$\begin{aligned}\beta &= P\{\text{error tipo II}\} = P\{\text{Aceptar equivocadamente } H_0\} \\ &= P\{\text{Aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\}\end{aligned}$$

Nota 66. Cuando H_0 es falsa y H_1 es verdadera, no siempre es posible especificar un valor exacto de μ , sino más bien un rango de posibles valores. En lugar de arriesgarse a tomar una decisión incorrecta, es mejor concluir que no hay evidencia suficiente para rechazar H_0 , es decir en lugar de aceptar H_0 , no rechazar H_0 .

La bondad de una prueba estadística se mide por el tamaño de α y β , ambas deben de ser pequeñas. Una manera muy efectiva de medir la potencia de la prueba es calculando el complemento del error tipo II:

$$\begin{aligned}1 - \beta &= P\{\text{Rechazar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{\text{Rechazar } H_0 \text{ cuando } H_1 \text{ es verdadera}\}\end{aligned}$$

Definición 53. La **potencia de la prueba**, $1 - \beta$, mide la capacidad de que la prueba funcione como se necesita.

47.14 Estimación por intervalos

Para la media

Recordemos que S^2 es un estimador insesgado de σ^2

Definición 54. Sean $\hat{\theta}_1$ y $\hat{\theta}_2$ dos estimadores insesgados de θ , parámetro poblacional. Si $\sigma_{\hat{\theta}_1}^2 < \sigma_{\hat{\theta}_2}^2$, decimos que $\hat{\theta}_1$ es un estimador más eficaz de θ que $\hat{\theta}_2$.

Algunas observaciones que es preciso realizar

- Para poblaciones normales, \bar{X} y \tilde{X} son estimadores insesgados de μ , pero con $\sigma_{\bar{X}}^2 < \sigma_{\tilde{X}}^2$.
- Para las estimaciones por intervalos de θ , un intervalo de la forma $\hat{\theta}_L < \theta < \hat{\theta}_U$, $\hat{\theta}_L$ y $\hat{\theta}_U$ dependen del valor de $\hat{\theta}$.
- Para $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$, si $n \rightarrow \infty$, entonces $\hat{\theta} \rightarrow \mu$.
- Para $\hat{\theta}$ se determinan $\hat{\theta}_L$ y $\hat{\theta}_U$ de modo tal que

$$P\{\hat{\theta}_L < \hat{\theta} < \hat{\theta}_U\} = 1 - \alpha, \quad (47.55)$$

con $\alpha \in (0, 1)$. Es decir, $\theta \in (\hat{\theta}_L, \hat{\theta}_U)$ es un intervalo de confianza del $100(1 - \alpha)\%$.

- De acuerdo con el TLC se espera que la distribución muestral de \bar{X} se distribuye aproximadamente normal con media $\mu_{\bar{X}} = \mu$ y desviación estándar $\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$.

Para $Z_{\alpha/2}$ se tiene $P\{-Z_{\alpha/2} < Z < Z_{\alpha/2}\} = 1 - \alpha$, donde $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$. Entonces $P\{-Z_{\alpha/2} < \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < Z_{\alpha/2}\} = 1 - \alpha$ es equivalente a $P\{\bar{X} - Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X} + Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\} = 1 - \alpha$

- f) Si \bar{X} es la media muestral de una muestra de tamaño n de una población con varianza conocida σ^2 , el intervalo de confianza de $100(1 - \alpha)\%$ para μ es $\mu \in \left(\bar{x} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right)$.
- g) Para muestras pequeñas de poblaciones no normales, no se puede esperar que el grado de confianza sea preciso.
- h) Para $n \geq 30$, con distribución de forma no muy sesgada, se pueden tener buenos resultados.

Teorema 5. Si \bar{X} es un estimador de μ , podemos tener $100(1 - \alpha)\%$ de confianza en que el error no excederá a $z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$, error entre \bar{X} y μ .

Teorema 6. Si \bar{X} es un estimador de μ , podemos tener $100(1 - \alpha)\%$ de confianza en que el error no excederá una cantidad e cuando el tamaño de la muestra es

$$n = \left(\frac{z_{\alpha/2}\sigma}{e}\right)^2.$$

Nota 67. Para intervalos unilaterales

$$P\left\{\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} < Z_{\alpha}\right\} = 1 - \alpha$$

equivalentemente

$$P\left\{\mu < \bar{X} + Z_{\alpha} \frac{\sigma}{\sqrt{n}}\right\} = 1 - \alpha.$$

Si \bar{X} es la media de una muestra aleatoria de tamaño n a partir de una población con varianza σ^2 , los límites de confianza unilaterales del $100(1 - \alpha)\%$ de confianza para μ están dados por

- Límite unilateral superior: $\bar{x} + z_{\alpha} \frac{\sigma}{\sqrt{n}}$
- Límite unilateral inferior: $\bar{x} - z_{\alpha} \frac{\sigma}{\sqrt{n}}$
- Para σ desconocida recordar que $T = \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t_{n-1}$, donde s es la desviación estándar de la muestra. Entonces

$$P\{-t_{\alpha/2} < T < t_{\alpha/2}\} = 1 - \alpha, \text{ equivalentemente}$$

$$P\left\{\bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}} < \mu < \bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}\right\} = 1 - \alpha.$$

- Un intervalo de confianza del $100(1 - \alpha)\%$ de confianza para μ , σ^2 desconocida y población normal es $\mu \in \left(\bar{x} - t_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}}\right)$, donde $t_{\alpha/2}$ es una t -student con $\nu = n - 1$ grados de libertad.
- Los límites unilaterales para μ con σ desconocida son $\bar{X} - t_{\alpha/2} \frac{s}{\sqrt{n}}$ y $\bar{X} + t_{\alpha/2} \frac{s}{\sqrt{n}}$.
- Cuando la población no es normal, σ desconocida y $n \geq 30$, σ se puede reemplazar por s para obtener el intervalo de confianza para muestras grandes:

$$\bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}.$$

- El estimador de \bar{X} de μ , σ desconocida, la varianza de $\sigma_{\bar{X}}^2 = \frac{\sigma^2}{n}$, el error estándar de \bar{X} es σ/\sqrt{n} .
- Si σ es desconocida y la población es normal, $s \rightarrow \sigma$ y se incluye el error estándar s/\sqrt{n} , entonces

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}.$$

Intervalos de confianza sobre la varianza

Supongamos que X se distribuye normal (μ, σ^2) , desconocidas. Sea X_1, X_2, \dots, X_n muestra aleatoria de tamaño n , s^2 la varianza muestral.

Se sabe que $X^2 = \frac{(n-1)s^2}{\sigma^2}$ se distribuye χ_{n-1}^2 grados de libertad. Su intervalo de confianza es

$$\begin{aligned} P \left\{ \chi_{1-\frac{\alpha}{2}, n-1}^2 \leq \chi^2 \leq \chi_{\frac{\alpha}{2}, n-1}^2 \right\} &= 1 - \alpha \\ P \left\{ \chi_{1-\frac{\alpha}{2}, n-1}^2 \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi_{\frac{\alpha}{2}, n-1}^2 \right\} &= 1 - \alpha \\ P \left\{ \frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right\} &= 1 - \alpha \end{aligned} \quad (47.56)$$

es decir

$$\sigma^2 \in \left[\frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2}, \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right] \quad (47.57)$$

los intervalos unilaterales son

$$\sigma^2 \in \left[\frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}, n-1}^2}, \infty \right) - \quad (47.58)$$

$$\sigma^2 \in \left(-\infty, \frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}, n-1}^2} \right] \quad (47.59)$$

Intervalos de confianza para proporciones

Supongamos que se tienen una muestra de tamaño n de una población grande pero finita, y supongamos que X , $X \leq n$, pertenecen a la clase de interés, entonces

$$\hat{p} = \frac{\bar{X}}{n}$$

es el estimador puntual de la proporción de la población que pertenece a dicha clase.

n y p son los parámetros de la distribución binomial, entonces $\hat{p} \sim N \left(p, \frac{p(1-p)}{n} \right)$ aproximadamente si p es distinto de 0 y 1; o si n es suficientemente grande. Entonces

$$Z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \sim N(0, 1), \text{ aproximadamente.}$$

entonces

$$\begin{aligned} 1 - \alpha &= P \left\{ -z_{\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \leq z_{\alpha/2} \right\} \\ &= P \left\{ \hat{p} - z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \right\} \end{aligned}$$

con $\sqrt{\frac{p(1-p)}{n}}$ error estándar del estimador puntual p . Una solución para determinar el intervalo de confianza del parámetro p (desconocido) es

$$1 - \alpha = P \left\{ \hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right\}$$

entonces los intervalos de confianza, tanto unilaterales como de dos colas son:

$$\begin{aligned} - p &\in \left(\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right) \\ - p &\in \left(-\infty, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right) \\ - p &\in \left(\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \infty \right) \end{aligned}$$

para minimizar el error estándar, se propone que el tamaño de la muestra sea $n = \left(\frac{z_{\alpha/2}}{E} \right)^2 p(1-p)$, donde $E = |p - \hat{p}|$.

47.15 Intervalos de confianza para dos muestras

Varianzas conocidas

Sean X_1 y X_2 variables aleatorias independientes. X_1 con media desconocida μ_1 y varianza conocida σ_1^2 ; y X_2 con media desconocida μ_2 y varianza conocida σ_2^2 . Se busca encontrar un intervalo de confianza de $100(1 - \alpha)\%$ de la diferencia entre medias μ_1 y μ_2 .

Sean $X_{11}, X_{12}, \dots, X_{1n_1}$ muestra aleatoria de n_1 observaciones de X_1 , y sean $X_{21}, X_{22}, \dots, X_{2n_2}$ muestra aleatoria de n_2 observaciones de X_2 .

Sean \bar{X}_1 y \bar{X}_2 , medias muestrales, entonces el estadístico

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0, 1), \quad (47.60)$$

si X_1 y X_2 son normales o aproximadamente normales si se aplican las condiciones del Teorema de Límite Central respectivamente.

Entonces se tiene

$$\begin{aligned} 1 - \alpha &= P \{ -Z_{\alpha/2} \leq Z \leq Z_{\alpha/2} \} \\ &= P \left\{ -Z_{\alpha/2} \leq \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \leq Z_{\alpha/2} \right\} \\ &= P \left\{ (\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \leq \mu_1 - \mu_2 \leq \right. \\ &\quad \left. (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right\} \end{aligned}$$

Entonces los intervalos de confianza unilaterales y de dos colas al $(1 - \alpha)\%$ de confianza son

$$\begin{aligned}
- \mu_1 - \mu_2 &\in \left[(\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right] \\
- \mu_1 - \mu_2 &\in \left[-\infty, (\bar{X}_1 - \bar{X}_2) + Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \right] \\
- \mu_1 - \mu_2 &\in \left[(\bar{X}_1 - \bar{X}_2) - Z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}, \infty \right]
\end{aligned}$$

Nota 68. Si σ_1 y σ_2 son conocidas, o por lo menos se conoce una aproximación, y los tamaños de las muestras n_1 y n_2 son iguales, $n_1 = n_2 = n$, se puede determinar el tamaño de la muestra para que el error al estimar $\mu_1 - \mu_2$ usando $\bar{X}_1 - \bar{X}_2$ sea menor que E (valor del error deseado) al $(1 - \alpha)\%$ de confianza. El tamaño n de la muestra requerido para cada muestra es

$$n = \left(\frac{Z_{\alpha/2}}{E} \right)^2 (\sigma_1^2 + \sigma_2^2).$$

Varianzas desconocidas

- Si $n_1, n_2 \geq 30$ se pueden utilizar los intervalos de la distribución normal para varianzas conocidas
- Si n_1, n_2 son muestras pequeñas, supongase que las poblaciones para X_1 y X_2 son normales con varianzas desconocidas y con base en el intervalo de confianza para distribuciones t -student

$$\sigma_1^2 = \sigma_2^2 = \sigma^2$$

Supongamos que X_1 es una variable aleatoria con media μ_1 y varianzas σ_1^2 , X_2 es una variable aleatoria con media μ_2 y varianzas σ_2^2 . Todos los parámetros son desconocidos. Sin embargo supóngase que es razonable considerar que $\sigma_1^2 = \sigma_2^2 = \sigma^2$.

Nuevamente sean X_1 y X_2 variables aleatorias independientes. X_1 con media desconocida μ_1 y varianzas muestral S_1^2 ; y X_2 con media desconocida μ_2 y varianzas muestral S_2^2 . Dado que S_1^2 y S_2^2 son estimadores de σ_1^2 , se propone el estimador S de σ^2 como

$$S_p^2 = \frac{(n_1 - 1) S_1^2 + (n_2 - 1) S_2^2}{n_1 + n_2 - 2},$$

entonces, el estadístico para $\mu_1 - \mu_2$ es

$$t_\nu = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

donde t_ν es una t de student con $\nu = n_1 + n_2 - 2$ grados de libertad.

Por lo tanto

$$\begin{aligned}
1 - \alpha &= P \{ -t_{\alpha/2, \nu} \leq t \leq t_{\alpha/2, \nu} \} \\
&= P \left\{ (\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \leq \right. \\
&\quad \left. t \leq (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right\}
\end{aligned}$$

luego, los intervalos de confianza del $(1 - \alpha)\%$ para $\mu_1 - \mu_2$ son

$$\begin{aligned}
- \mu_1 - \mu_2 &\in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right] \\
- \mu_1 - \mu_2 &\in \left[-\infty, (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} \right] \\
- \mu_1 - \mu_2 &\in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}, \infty \right]
\end{aligned}$$

$$\sigma_1^2 \neq \sigma_2^2$$

Si no se tiene certeza de que $\sigma_1^2 = \sigma_2^2$, se propone el estadístico

$$t^* = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \quad (47.61)$$

que se distribuye t -student con ν grados de libertad, donde

$$\nu = \frac{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^2}{\frac{S_1^2/n_1}{n_1+1} + \frac{S_2^2/n_2}{n_2+1}} - 2$$

Entonces el intervalo de confianza de aproximadamente el $100(1 - \alpha)\%$ para $\mu_1 - \mu_2$ con $\sigma_1^2 \neq \sigma_2^2$ es

$$\begin{aligned}
\mu_1 - \mu_2 &\in \left[(\bar{X}_1 - \bar{X}_2) - t_{\alpha/2, \nu} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}, \right. \\
&\quad \left. (\bar{X}_1 - \bar{X}_2) + t_{\alpha/2, \nu} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}} \right]
\end{aligned}$$

47.16 Intervalos de confianza para razón de Varianzas

Supongamos que se toman dos muestras aleatorias independientes de las dos poblaciones de interés.

Sean X_1 y X_2 variables normales independientes con medias desconocidas μ_1 y μ_2 y varianzas desconocidas σ_1^2 y σ_2^2 respectivamente. Se busca un intervalo de confianza de $100(1 - \alpha)\%$ para σ_1^2/σ_2^2 . Supongamos n_1 y n_2 muestras aleatorias de X_1 y X_2 y sean S_1^2 y S_2^2 varianzas muestrales. Se sabe que

$$F = \frac{S_2^2/\sigma_2^2}{S_1^2/\sigma_1^2}$$

se distribuye F con $n_2 - 1$ y $n_1 - 1$ grados de libertad.

Por lo tanto

$$\begin{aligned}
P \left\{ F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq F \leq F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right\} &= 1 - \alpha \\
P \left\{ F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq \frac{S_2^2/\sigma_2^2}{S_1^2/\sigma_1^2} \leq F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right\} &= 1 - \alpha
\end{aligned}$$

por lo tanto

$$P \left\{ \frac{S_1^2}{S_2^2} F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} \leq \frac{\sigma_1^2}{\sigma_2^2} \leq \frac{S_1^2}{S_2^2} F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right\} = 1 - \alpha$$

entonces

$$\frac{\sigma_1^2}{\sigma_2^2} \in \left[\frac{S_1^2}{S_2^2} F_{1-\frac{\alpha}{2}, n_2-1, n_1-1}, \frac{S_1^2}{S_2^2} F_{\frac{\alpha}{2}, n_2-1, n_1-1} \right]$$

donde

$$F_{1-\frac{\alpha}{2}, n_2-1, n_1-1} = \frac{1}{F_{\frac{\alpha}{2}, n_2-1, n_1-1}}$$

47.17 Intervalos de confianza para diferencia de proporciones

Sean dos proporciones de interés p_1 y p_2 . Se busca un intervalo para $p_1 - p_2$ al $100(1 - \alpha)\%$.

Sean dos muestras independientes de tamaño n_1 y n_2 de poblaciones infinitas de modo que X_1 y X_2 variables aleatorias binomiales independientes con parámetros (n_1, p_1) y (n_2, p_2) .

X_1 y X_2 son el número de observaciones que pertenecen a la clase de interés correspondientes. Entonces $\hat{p}_1 = \frac{X_1}{n_1}$ y $\hat{p}_2 = \frac{X_2}{n_2}$ son estimadores de p_1 y p_2 respectivamente. Supongamos que se cumple la aproximación normal a la binomial, entonces

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{n_1} - \frac{p_2(1-p_2)}{n_2}}} \sim N(0, 1) \text{ aproximadamente}$$

entonces

$$\begin{aligned} (\hat{p}_1 - \hat{p}_2) - Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} &\leq p_1 - p_2 \\ &\leq (\hat{p}_1 - \hat{p}_2) + Z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} - \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \end{aligned}$$

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado
- el interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.
- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

47.18 2. Pruebas de Hipótesis

47.18.1 2.1 Tipos de errores

Una prueba de hipótesis está formada por cinco partes

- La hipótesis nula, denotada por H_0 .
- La hipótesis alterativa, denorada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 55. Las dos hipótesis en competencias son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones:

- Rechazar H_0 y concluir que H_1 es verdadera.
- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 28. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 29. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 56. El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadística de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**.

Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**.

Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente.

Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 57. Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

47.19 2.2 Muestras grandes: una media poblacional

47.19.1 2.2.1 Cálculo de valor p

Definición 58. El **valor de p** (*p-value*) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.

Nota 69. Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.

Nota 70. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 59. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados

p	H_0	Significativa
$p < 0.01$	Rechazar	Altamente
$0.01 \leq p < 0.05$	Rechazar	Estadísticamente
$0.05 \leq p < 0.1$	No rechazar	Tendencia estadística
$0.1 \leq p$	No rechazar	No son estadísticamente

Nota 71. Para determinar el valor de p , encontrar el área en la cola después del estadístico de prueba. Si la prueba es de una cola, este es el valor de p . Si es de dos colas, éste valor encontrado es la mitad del valor de p . Rechazar H_0 cuando el valor de $p < \alpha$.

Hay dos tipos de errores al realizar una prueba de hipótesis

	H_0 es Verdadera	H_0 es Falsa
Rechazar H_0	Error tipo I	✓
Aceptar H_0	✓	Error tipo II

Definición 60. La probabilidad de cometer el error tipo II se define por β donde

$$\begin{aligned}\beta &= P\{\text{error tipo II}\} = P\{\text{Aceptar equivocadamente } H_0\} \\ &= P\{\text{Aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\}\end{aligned}$$

Nota 72. Cuando H_0 es falsa y H_1 es verdadera, no siempre es posible especificar un valor exacto de μ , sino más bien un rango de posibles valores. En lugar de arriesgarse a tomar una decisión incorrecta, es mejor concluir que no hay evidencia suficiente para rechazar H_0 , es decir en lugar de aceptar H_0 , no rechazar H_0 .

La bondad de una prueba estadística se mide por el tamaño de α y β , ambas deben de ser pequeñas. Una manera muy efectiva de medir la potencia de la prueba es calculando el complemento del error tipo II:

$$\begin{aligned}1 - \beta &= P\{\text{Rechazar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{\text{Rechazar } H_0 \text{ cuando } H_1 \text{ es verdadera}\}\end{aligned}$$

Definición 61. La **potencia de la prueba**, $1 - \beta$, mide la capacidad de que la prueba funciona como se necesita.

Ejemplo 30. La producción diaria de una planta química local ha promediado 880 toneladas en los últimos años. A la gerente de control de calidad le gustaría saber si este promedio ha cambiado en meses recientes. Ella selecciona al azar 50 días de la base de datos computarizada y calcula el promedio y la desviación estándar de las $n = 50$ producciones como $\bar{x} = 871$ toneladas y $s = 21$ toneladas, respectivamente. Pruebe la hipótesis apropiada usando $\alpha = 0.05$.

Solución 6. La hipótesis nula apropiada es:

$$\begin{aligned}H_0 &: \mu = 880 \\ &\quad \text{y la hipótesis alternativa } H_1 \text{ es} \\ H_1 &: \mu \neq 880\end{aligned}$$

el estimador puntual para μ es \bar{x} , entonces el estadístico de prueba es

$$\begin{aligned}z &= \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \\ &= \frac{871 - 880}{21/\sqrt{50}} = -3.03\end{aligned}$$

Solución 7. Para esta prueba de dos colas, hay que determinar los dos valores de $z_{\alpha/2}$, es decir, $z_{\alpha/2} = \pm 1.96$, como $z > z_{\alpha/2}$, z cae en la zona de rechazo, por lo tanto la gerente puede rechazar la hipótesis nula y concluir que el promedio efectivamente ha cambiado. La probabilidad de rechazar H_0 cuando esta es verdadera es de 0.05.

Recordemos que el valor observado del estadístico de prueba es $z = -3.03$, la región de rechazo más pequeña que puede usarse y todavía seguir rechazando H_0 es $|z| > 3.03$, entonces $p = 2(0.012) = 0.0024$, que a su vez es menor que el nivel de significancia α asignado inicialmente, y además los resultados son **altamente significativos**.

Finalmente determinemos la potencia de la prueba cuando μ en realidad es igual a 870 toneladas.

Recordar que la región de aceptación está entre -1.96 y 1.96 , para $\mu = 880$, equivalentemente

$$874.18 < \bar{x} < 885.82$$

β es la probabilidad de aceptar H_0 cuando $\mu = 870$, calculemos los valores de z correspondientes a 874.18 y 885.82 Entonces

$$z_1 = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{874.18 - 870}{21/\sqrt{50}} = 1.41$$

$$z_1 = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{885.82 - 870}{21/\sqrt{50}} = 5.33$$

por lo tanto

$$\begin{aligned}\beta &= P\{\text{aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{874.18 < \mu < 885.82 \text{ cuando } \mu = 870\} \\ &= P\{1.41 < z < 5.33\} = P\{1.41 < z\} \\ &= 1 - 0.9207 = 0.0793\end{aligned}$$

entonces, la potencia de la prueba es

$$1 - \beta = 1 - 0.0793 = 0.9207$$

que es la probabilidad de rechazar correctamente H_0 cuando H_0 es falsa.

Determinar la potencia de la prueba para distintos valores de H_1 y graficarlos, *curva de potencia*

H_1	$(1 - \beta)$
865	
870	
872	
875	
877	
880	
883	
885	
888	
890	
895	

- Encontrar las regiones de rechazo para el estadístico z , para una prueba de
 - dos colas para $\alpha = 0.01, 0.05, 0.1$
 - una cola superior para $\alpha = 0.01, 0.05, 0.1$
 - una cola inferior para $\alpha = 0.01, 0.05, 0.1$
- Suponga que el valor del estadístico de prueba es
 - $z = -2.41$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = 2.16$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = 1.15$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = -2.78$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = -1.81$, sacar las conclusiones correspondientes para los incisos anteriores.
- Encuentre el valor de p para las pruebas de hipótesis correspondientes a los valores de z del ejercicio anterior.
- Para las pruebas dadas en el ejercicio 2, utilice el valor de p , determinado en el ejercicio 3, para determinar la significancia de los resultados.

5. Una muestra aleatoria de $n = 45$ observaciones de una población con media $\bar{x} = 2.4$, y desviación estándar $s = 0.29$. Suponga que el objetivo es demostrar que la media poblacional μ excede 2.3.
 - a) Defina la hipótesis nula y alternativa para la prueba.
 - b) Determine la región de rechazo para un nivel de significancia de: $\alpha = 0.1, 0.05, 0.01$.
 - c) Determine el error estándar de la media muestral.
 - d) Calcule el valor de p para los estadísticos de prueba definidos en los incisos anteriores.
 - e) Utilice el valor de p para sacar una conclusión al nivel de significancia α .
 - f) Determine el valor de β cuando $\mu = 2.5$
 - g) Graficar la curva de potencia para la prueba.

47.19.2 2.2.2 Prueba de hipótesis para la diferencia entre dos medias poblacionales

El estadístico que resume la información muestral respecto a la diferencia en medias poblacionales $(\mu_1 - \mu_2)$ es la diferencia de las medias muestrales $(\bar{x}_1 - \bar{x}_2)$, por tanto al probar la diferencia entre las medias muestrales se verifica que la diferencia real entre las medias poblacionales difiere de un valor especificado, $(\mu_1 - \mu_2) = D_0$, se puede usar el error estándar de $(\bar{x}_1 - \bar{x}_2)$, es decir

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

cuyo estimador está dado por

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

El procedimiento para muestras grandes es:

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

- | | Prueba de una Cola | Prueba de dos colas |
|---------------------------------|--|----------------------------------|
| 2) Hipótesis Alternativa | $H_1 : (\mu_1 - \mu_2) > D_0$
$H_1 : (\mu_1 - \mu_2) < D_0$ | $H_1 : (\mu_1 - \mu_2) \neq D_0$ |

- 3) Estadístico de prueba:

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

- | | Prueba de una Cola | Prueba de dos colas |
|---|--|--|
| 4) Región de rechazo: rechazar H_0 cuando | $z > z_0$
$z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$
cuando $p < \alpha$ | $z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$ |

Ejemplo 31. Para determinar si ser propietario de un automóvil afecta el rendimiento académico de un estudiante, se tomaron dos muestras aleatorias de 100 estudiantes varones. El promedio de calificaciones para los $n_1 = 100$ no propietarios de un auto tuvieron un promedio y varianza de $\bar{x}_1 = 2.7$ y $s_1^2 = 0.36$, respectivamente, mientras que para la segunda muestra con $n_2 = 100$ propietarios de un auto, se tiene $\bar{x}_2 = 2.54$ y $s_2^2 = 0.4$. Los datos presentan suficiente evidencia para indicar una diferencia en la media en el rendimiento académico entre propietarios y no propietarios de un automóvil? Hacer pruebas para $\alpha = 0.01, 0.05$ y $\alpha = 0.1$.

Solución 8. – Solución utilizando la técnica de regiones de rechazo: realizando las operaciones $z = 1.84$, determinar si excede los valores de $z_{\alpha/2}$.

– Solución utilizando el p -value: Calcular el valor de p , la probabilidad de que z sea mayor que $z = 1.84$ o menor que $z = -1.84$, se tiene que $p = 0.0658$. Concluir.

- Si el intervalo de confianza que se construye contiene el valor del parámetro especificado por H_0 , entonces ese valor es uno de los posibles valores del parámetro y H_0 no debe ser rechazada.
- Si el valor hipotético se encuentra fuera de los límites de confianza, la hipótesis nula es rechazada al nivel de significancia α .

1. Del libro Mendenhall resolver los ejercicios 9.18, 9.19 y 9.20(Mendenhall).
2. Del libro Mendenhall resolver los ejercicios: 9.23, 9.26 y 9.28.

47.19.3 2.2.3 Prueba de Hipótesis para una Proporción Binomial

Para una muestra aleatoria de n intentos idénticos, de una población binomial, la proporción muestral \hat{p} tiene una distribución aproximadamente normal cuando n es grande, con media p y error estándar

$$SE = \sqrt{\frac{pq}{n}}.$$

La prueba de hipótesis de la forma

$$H_0 : p = p_0$$

$$H_1 : p > p_0, \text{ o } p < p_0 \text{ o } p \neq p_0$$

El estadístico de prueba se construye con el mejor estimador de la proporción verdadera, \hat{p} , con el estadístico de prueba z , que se distribuye normal estándar.

El procedimiento es

- 1) Hipótesis nula: $H_0 : p = p_0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis alternativa	$H_1 : p > p_0$ $H_1 : p < p_0$	$p \neq p_0$

- 3) Estadístico de prueba:

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{pq}{n}}}, \hat{p} = \frac{x}{n}$$

donde x es el número de éxitos en n intentos binomiales.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_0$ $z < -z_\alpha$ cuando $H_1 : p < p_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Ejemplo 32. A cualquier edad, alrededor del 20% de los adultos de cierto país realiza actividades de acondicionamiento físico al menos dos veces por semana. En una encuesta local de $n = 100$ adultos de más de 40 a nos, un total de 15 personas indicaron que realizaron actividad física al menos dos veces por semana. Estos datos indican que el porcentaje de participación para adultos de más de 40 a nos de edad es considerablemente menor a la cifra del 20%? Calcule el valor de p y úselo para sacar las conclusiones apropiadas.

1. Resolver los ejercicios: 9.30, 9.32, 9.33, 9.35 y 9.39.

47.19.4 2.2.4 Prueba de Hipótesis diferencia entre dos Proporciones Binomiales

Nota 73. Cuando se tienen dos muestras aleatorias independientes de dos poblaciones binomiales, el objetivo del experimento puede ser la diferencia $(p_1 - p_2)$ en las proporciones de individuos u objetos que poseen una característica específica en las dos poblaciones. En este caso se pueden utilizar los estimadores de las dos proporciones $(\hat{p}_1 - \hat{p}_2)$ con error estándar dado por

$$SE = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

considerando el estadístico z con un nivel de significancia $(1 - \alpha) 100\%$

Nota 74. La hipótesis nula a probarse es de la forma

$H_0: p_1 = p_2$ o equivalentemente $(p_1 - p_2) = 0$, contra una hipótesis alternativa H_1 de una o dos colas.

Nota 75. Para estimar el error estándar del estadístico z , se debe de utilizar el hecho de que suponiendo que H_0 es verdadera, las dos proporciones son iguales a algún valor común, p . Para obtener el mejor estimador de p es

$$p = \frac{\text{número total de éxitos}}{\text{Número total de pruebas}} = \frac{x_1 + x_2}{n_1 + n_2}$$

1) **Hipótesis Nula:** $H_0 : (p_1 - p_2) = 0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 :$	$H_1 : (p_1 - p_2) > 0$ $H_1 : (p_1 - p_2) < 0$	$H_1 : (p_1 - p_2) \neq 0$

3) Estadístico de prueba:

$$z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{pq}{n_1} + \frac{pq}{n_2}}}$$

donde $\hat{p}_1 = x_1/n_1$ y $\hat{p}_2 = x_2/n_2$, dado que el valor común para p_1 y p_2 es p , entonces $\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$ y por tanto el estadístico de prueba es

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_\alpha$ $z < -z_\alpha$ cuando $H_1 : p < p_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Ejemplo 33. Los registros de un hospital, indican que 52 hombres de una muestra de 1000 contra 23 mujeres de una muestra de 1000 fueron ingresados por enfermedad del corazón. Estos datos presentan suficiente evidencia para indicar un porcentaje más alto de enfermedades del corazón entre hombres ingresados al hospital?, utilizar distintos niveles de confianza de α .

1. Resolver los ejercicios 9.42
2. Resolver los ejercicios: 9.45, 9.48, 9.50

47.20 2.3 Muestras Pequeñas

47.20.1 2.3.1 Una media poblacional

1) **Hipótesis Nula:** $H_0 : \mu = \mu_0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 :$	$H_1 : \mu > \mu_0$ $H_1 : \mu < \mu_0$	$H_1 : \mu \neq \mu_0$

3) Estadístico de prueba:

$$t = \frac{\bar{x} - \mu_0}{\sqrt{\frac{s^2}{n}}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Ejemplo 34. Las etiquetas en latas de un galón de pintura por lo general indican el tiempo de secado y el área puede cubrir una capa. Casi todas las marcas de pintura indican que, en una capa, un galón cubrirá entre 250 y 500 pies cuadrados, dependiendo de la textura de la superficie a pintarse, un fabricante, sin embargo afirma que un galón de su pintura cubrirá 400 pies cuadrados de área superficial. Para probar su afirmación, una muestra aleatoria de 10 latas de un galón de pintura blanca se empleó para pintar 10 áreas idénticas usando la misma clase de equipo. Las áreas reales en pies cuadrados cubiertas por estos 10 galones de pintura se dan a continuación:

310	311	412	368	447
376	303	410	365	350

Ejemplo 35. Los datos presentan suficiente evidencia para indicar que el promedio de la cobertura difiere de 400 pies cuadrados? encuentre el valor de p para la prueba y úselo para evaluar la significancia de los resultados.

1. Resolver los ejercicios: 10.2, 10.3, 10.5, 10.7, 10.9, 10.13 y 10.16

47.20.2 2.3.2 Diferencia entre dos medias poblacionales: M.A.I.

Nota 76. Cuando los tamaños de muestra son pequeños, no se puede asegurar que las medias muestrales sean normales, pero si las poblaciones originales son normales, entonces la distribución muestral de la diferencia de las medias muestrales, $(\bar{x}_1 - \bar{x}_2)$, será normal con media $(\mu_1 - \mu_2)$ y error estándar

$$ES = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\mu_1 - \mu_2) > D_0$ $H_1 : (\mu_1 - \mu_2) < D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$

- 3) Estadístico de prueba:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

donde

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

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_0$ $z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

47.20.3 2.3.3 Diferencia entre dos medias poblacionales: Diferencias Pareadas

- 1) **Hipótesis Nula:** $H_0 : \mu_d = 0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 : \mu_d$	$H_1 : \mu_d > 0$ $H_1 : \mu_d < 0$	$H_1 : \mu_d \neq 0$

- 3) Estadístico de prueba:

$$t = \frac{\bar{d}}{\sqrt{\frac{s_d^2}{n}}}$$

donde n es el número de diferencias pareadas, \bar{d} es la media de las diferencias muestrales, y s_d es la desviación estándar de las diferencias muestrales.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

47.20.4 2.3.4 Inferencias con respecto a la Varianza Poblacional

- 1) **Hipótesis Nula:** $H_0 : \sigma^2 = \sigma_0^2$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: H_1	$H_1 : \sigma^2 > \sigma_0^2$ $H_1 : \sigma^2 < \sigma_0^2$	$H_1 : \sigma^2 \neq \sigma_0^2$

- 3) Estadístico de prueba:

$$\chi^2 = \frac{(n - 1)s^2}{\sigma_0^2}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$\chi^2 > \chi_\alpha^2$ $\chi^2 < \chi_{(1-\alpha)}^2$ cuando $H_1 : \chi^2 < \chi_0^2$ cuando $p < \alpha$	$\chi^2 > \chi_{\alpha/2}^2$ o $\chi^2 < \chi_{(1-\alpha/2)}^2$

Los valores críticos de χ^2 , están basados en $(n_1 +)$ grados de libertad.

47.20.5 2.3.5 Comparación de dos varianzas poblacionales

- 1) **Hipótesis Nula** $H_0 : (\sigma_1^2 - \sigma_2^2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\sigma_1^2 - \sigma_2^2) > D_0$ $H_1 : (\sigma_1^2 - \sigma_2^2) < D_0$	$H_1 : (\sigma_1^2 - \sigma_2^2) \neq D_0$

- 3) Estadístico de prueba:

$$F = \frac{s_1^2}{s_2^2}$$

donde s_1^2 es la varianza muestral más grande.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$F > F_\alpha$ cuando $p < \alpha$	$F > F_{\alpha/2}$

47.21 2. Pruebas de Hipótesis

47.21.1 2.1 Tipos de errores

- Una hipótesis estadística es una afirmación acerca de la distribución de probabilidad de una variable aleatoria, a menudo involucran uno o más parámetros de la distribución.
- Las hipótesis son afirmaciones respecto a la población o distribución bajo estudio, no en torno a la muestra.
- La mayoría de las veces, la prueba de hipótesis consiste en determinar si la situación experimental ha cambiado
- el interés principal es decidir sobre la veracidad o falsedad de una hipótesis, a este procedimiento se le llama *prueba de hipótesis*.
- Si la información es consistente con la hipótesis, se concluye que esta es verdadera, de lo contrario que con base en la información, es falsa.

Una prueba de hipótesis está formada por cinco partes

- La hipótesis nula, denotada por H_0 .
- La hipótesis alterativa, denorada por H_1 .
- El estadístico de prueba y su valor p .
- La región de rechazo.
- La conclusión.

Definición 62. Las dos hipótesis en competencias son la **hipótesis alternativa** H_1 , usualmente la que se desea apoyar, y la **hipótesis nula** H_0 , opuesta a H_1 .

En general, es más fácil presentar evidencia de que H_1 es cierta, que demostrar que H_0 es falsa, es por eso que por lo regular se comienza suponiendo que H_0 es cierta, luego se utilizan los datos de la muestra para decidir si existe evidencia a favor de H_1 , más que a favor de H_0 , así se tienen dos conclusiones:

- Rechazar H_0 y concluir que H_1 es verdadera.
- Aceptar, no rechazar, H_0 como verdadera.

Ejemplo 36. Se desea demostrar que el salario promedio por hora en cierto lugar es distinto de 19usd, que es el promedio nacional. Entonces $H_1 : \mu \neq 19$, y $H_0 : \mu = 19$.

A esta se le denomina **Prueba de hipótesis de dos colas**.

Ejemplo 37. Un determinado proceso produce un promedio de 5% de piezas defectuosas. Se está interesado en demostrar que un simple ajuste en una máquina reducirá p , la proporción de piezas defectuosas producidas en este proceso. Entonces se tiene $H_0 : p < 0.3$ y $H_1 : p = 0.03$. Si se puede rechazar H_0 , se concluye que el proceso ajustado produce menos del 5% de piezas defectuosas.

A esta se le denomina **Prueba de hipótesis de una cola**.

La decisión de rechazar o aceptar la hipótesis nula está basada en la información contenida en una muestra proveniente de la población de interés. Esta información tiene estas formas

- **Estadístico de prueba:** un sólo número calculado a partir de la muestra.
- **p -value:** probabilidad calculada a partir del estadístico de prueba.

Definición 63. El p -value es la probabilidad de observar un estadístico de prueba tanto o más alejado del valor observado, si en realidad H_0 es verdadera. Valores grandes del estadística de prueba y valores pequeños de p significan que se ha observado un evento muy poco probable, si H_0 en realidad es verdadera.

Todo el conjunto de valores que puede tomar el estadístico de prueba se divide en dos regiones. Un conjunto, formado de valores que apoyan la hipótesis alternativa y llevan a rechazar H_0 , se denomina **región de rechazo**. El otro, conformado por los valores que sustentan la hipótesis nula, se le denomina **región de aceptación**.

Cuando la región de rechazo está en la cola izquierda de la distribución, la prueba se denomina **prueba lateral izquierda**. Una prueba con región de rechazo en la cola derecha se le llama **prueba lateral derecha**.

Si el estadístico de prueba cae en la región de rechazo, entonces se rechaza H_0 . Si el estadístico de prueba cae en la región de aceptación, entonces la hipótesis nula se acepta o la prueba se juzga como no concluyente.

Dependiendo del nivel de confianza que se desea agregar a las conclusiones de la prueba, y el **nivel de significancia** α , el riesgo que está dispuesto a correr si se toma una decisión incorrecta.

Definición 64. Un **error de tipo I** para una prueba estadística es el error que se tiene al rechazar la hipótesis nula cuando es verdadera. El **nivel de significancia** para una prueba estadística de hipótesis es

$$\begin{aligned}\alpha &= P\{\text{error tipo I}\} = P\{\text{rechazar equivocadamente } H_0\} \\ &= P\{\text{rechazar } H_0 \text{ cuando } H_0 \text{ es verdadera}\}\end{aligned}$$

Este valor α representa el valor máximo de riesgo tolerable de rechazar incorrectamente H_0 . Una vez establecido el nivel de significancia, la región de rechazo se define para poder determinar si se rechaza H_0 con un cierto nivel de confianza.

47.22 2.2 Muestras grandes: una media poblacional

47.22.1 2.2.1 Cálculo de valor p

Definición 65. El **valor de p** (*p-value*) o nivel de significancia observado de un estadístico de prueba es el valor más pequeño de α para el cual H_0 se puede rechazar. El riesgo de cometer un error tipo I, si H_0 es rechazada con base en la información que proporciona la muestra.

Nota 77. Valores pequeños de p indican que el valor observado del estadístico de prueba se encuentra alejado del valor hipotético de μ , es decir se tiene evidencia de que H_0 es falsa y por tanto debe de rechazarse.

Nota 78. Valores grandes de p indican que el estadístico de prueba observado no está alejado de la media hipotética y no apoya el rechazo de H_0 .

Definición 66. Si el valor de p es menor o igual que el nivel de significancia α , determinado previamente, entonces H_0 es rechazada y se puede concluir que los resultados son estadísticamente significativos con un nivel de confianza del $100(1 - \alpha)\%$.

Es usual utilizar la siguiente clasificación de resultados

p	H_0	Significativa
$p < 0.01$	Rechazar	Altamente
$0.01 \leq p < 0.05$	Rechazar	Estadísticamente
$0.05 \leq p < 0.1$	No rechazar	Tendencia estadística
$0.1 \leq p$	No rechazar	No son estadísticamente

Nota 79. Para determinar el valor de p , encontrar el área en la cola después del estadístico de prueba. Si la prueba es de una cola, este es el valor de p . Si es de dos colas, éste valor encontrado es la mitad del valor de p . Rechazar H_0 cuando el valor de $p < \alpha$.

Hay dos tipos de errores al realizar una prueba de hipótesis

	H_0 es Verdadera	H_0 es Falsa
Rechazar H_0	Error tipo I	✓
Aceptar H_0	✓	Error tipo II

Definición 67. La probabilidad de cometer el error tipo II se define por β donde

$$\begin{aligned}\beta &= P\{\text{error tipo II}\} = P\{\text{Aceptar equivocadamente } H_0\} \\ &= P\{\text{Aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\}\end{aligned}$$

Nota 80. Cuando H_0 es falsa y H_1 es verdadera, no siempre es posible especificar un valor exacto de μ , sino más bien un rango de posibles valores. En lugar de arriesgarse a tomar una decisión incorrecta, es mejor concluir que no hay evidencia suficiente para rechazar H_0 , es decir en lugar de aceptar H_0 , no rechazar H_0 .

La bondad de una prueba estadística se mide por el tamaño de α y β , ambas deben de ser pequeñas. Una manera muy efectiva de medir la potencia de la prueba es calculando el complemento del error tipo II:

$$\begin{aligned}1 - \beta &= P\{\text{Rechazar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{\text{Rechazar } H_0 \text{ cuando } H_1 \text{ es verdadera}\}\end{aligned}$$

Definición 68. La **potencia de la prueba**, $1 - \beta$, mide la capacidad de que la prueba funciona como se necesita.

Ejemplo 38. La producción diaria de una planta química local ha promediado 880 toneladas en los últimos años. A la gerente de control de calidad le gustaría saber si este promedio ha cambiado en meses recientes. Ella selecciona al azar 50 días de la base de datos computarizada y calcula el promedio y la desviación estándar de las $n = 50$ producciones como $\bar{x} = 871$ toneladas y $s = 21$ toneladas, respectivamente. Pruebe la hipótesis apropiada usando $\alpha = 0.05$.

Solución 9. La hipótesis nula apropiada es:

$$\begin{aligned} H_0 &: \mu = 880 \\ &\text{y la hipótesis alternativa } H_1 \text{ es} \\ H_1 &: \mu \neq 880 \end{aligned}$$

el estimador puntual para μ es \bar{x} , entonces el estadístico de prueba es

$$\begin{aligned} z &= \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \\ &= \frac{871 - 880}{21/\sqrt{50}} = -3.03 \end{aligned}$$

Solución 10. Para esta prueba de dos colas, hay que determinar los dos valores de $z_{\alpha/2}$, es decir, $z_{\alpha/2} = \pm 1.96$, como $z > z_{\alpha/2}$, z cae en la zona de rechazo, por lo tanto la gerente puede rechazar la hipótesis nula y concluir que el promedio efectivamente ha cambiado. La probabilidad de rechazar H_0 cuando esta es verdadera es de 0.05.

Recordemos que el valor observado del estadístico de prueba es $z = -3.03$, la región de rechazo más pequeña que puede usarse y todavía seguir rechazando H_0 es $|z| > 3.03$, entonces $p = 2(0.012) = 0.0024$, que a su vez es menor que el nivel de significancia α asignado inicialmente, y además los resultados son **altamente significativos**.

Finalmente determinemos la potencia de la prueba cuando μ en realidad es igual a 870 toneladas.

Recordar que la región de aceptación está entre -1.96 y 1.96 , para $\mu = 880$, equivalentemente

$$874.18 < \bar{x} < 885.82$$

β es la probabilidad de aceptar H_0 cuando $\mu = 870$, calculemos los valores de z correspondientes a 874.18 y 885.82. Entonces

$$\begin{aligned} z_1 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{874.18 - 870}{21/\sqrt{50}} = 1.41 \\ z_1 &= \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{885.82 - 870}{21/\sqrt{50}} = 5.33 \end{aligned}$$

por lo tanto

$$\begin{aligned} \beta &= P\{\text{aceptar } H_0 \text{ cuando } H_0 \text{ es falsa}\} \\ &= P\{874.18 < \mu < 885.82 \text{ cuando } \mu = 870\} \\ &= P\{1.41 < z < 5.33\} = P\{1.41 < z\} \\ &= 1 - 0.9207 = 0.0793 \end{aligned}$$

entonces, la potencia de la prueba es

$$1 - \beta = 1 - 0.0793 = 0.9207$$

que es la probabilidad de rechazar correctamente H_0 cuando H_0 es falsa.

Determinar la potencia de la prueba para distintos valores de H_1 y graficarlos, *curva de potencia*

H_1	$(1 - \beta)$
865	
870	
872	
875	
877	
880	
883	
885	
888	
890	
895	

- Encontrar las regiones de rechazo para el estadístico z , para una prueba de
 - dos colas para $\alpha = 0.01, 0.05, 0.1$
 - una cola superior para $\alpha = 0.01, 0.05, 0.1$
 - una cola inferior para $\alpha = 0.01, 0.05, 0.1$
- Suponga que el valor del estadístico de prueba es
 - $z = -2.41$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = 2.16$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = 1.15$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = -2.78$, sacar las conclusiones correspondientes para los incisos anteriores.
 - $z = -1.81$, sacar las conclusiones correspondientes para los incisos anteriores.
- Encuentre el valor de p para las pruebas de hipótesis correspondientes a los valores de z del ejercicio anterior.
- Para las pruebas dadas en el ejercicio 2, utilice el valor de p , determinado en el ejercicio 3, para determinar la significancia de los resultados.
- Una muestra aleatoria de $n = 45$ observaciones de una población con media $\bar{x} = 2.4$, y desviación estándar $s = 0.29$. Suponga que el objetivo es demostrar que la media poblacional μ excede 2.3.
 - Defina la hipótesis nula y alternativa para la prueba.
 - Determine la región de rechazo para un nivel de significancia de: $\alpha = 0.1, 0.05, 0.01$.
 - Determine el error estándar de la media muestral.
 - Calcule el valor de p para los estadísticos de prueba definidos en los incisos anteriores.
 - Utilice el valor de p para sacar una conclusión al nivel de significancia α .
 - Determine el valor de β cuando $\mu = 2.5$
 - Graficar la curva de potencia para la prueba.

47.22.2 2.2.2 Prueba de hipótesis para la diferencia entre dos medias poblacionales

El estadístico que resume la información muestral respecto a la diferencia en medias poblacionales $(\mu_1 - \mu_2)$ es la diferencia de las medias muestrales $(\bar{x}_1 - \bar{x}_2)$, por tanto al probar la diferencia entre las medias muestrales se verifica que la diferencia real entre las medias poblacionales difiere de un valor especificado, $(\mu_1 - \mu_2) = D_0$, se puede usar el error estándar de $(\bar{x}_1 - \bar{x}_2)$, es decir

$$\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

cuyo estimador está dado por

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

El procedimiento para muestras grandes es:

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

- 2) **Hipótesis Alternativa**

Prueba de una Cola	Prueba de dos colas
$H_1 : (\mu_1 - \mu_2) > D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$
$H_1 : (\mu_1 - \mu_2) < D_0$	

- 3) Estadístico de prueba:

$$z = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

- 4) Región de rechazo: rechazar H_0 cuando

Prueba de una Cola	Prueba de dos colas
$z > z_0$	
$z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$
cuando $p < \alpha$	

Ejemplo 39. Para determinar si ser propietario de un automóvil afecta el rendimiento académico de un estudiante, se tomaron dos muestras aleatorias de 100 estudiantes varones. El promedio de calificaciones para los $n_1 = 100$ no propietarios de un auto tuvieron un promedio y varianza de $\bar{x}_1 = 2.7$ y $s_1^2 = 0.36$, respectivamente, mientras que para la segunda muestra con $n_2 = 100$ propietarios de un auto, se tiene $\bar{x}_2 = 2.54$ y $s_2^2 = 0.4$. Los datos presentan suficiente evidencia para indicar una diferencia en la media en el rendimiento académico entre propietarios y no propietarios de un automóvil? Hacer pruebas para $\alpha = 0.01, 0.05$ y $\alpha = 0.1$.

Solución 11. – Solución utilizando la técnica de regiones de rechazo: realizando las operaciones $z = 1.84$, determinar si excede los valores de $z_{\alpha/2}$.

- Solución utilizando el p -value: Calcular el valor de p , la probabilidad de que z sea mayor que $z = 1.84$ o menor que $z = -1.84$, se tiene que $p = 0.0658$. Concluir.
- Si el intervalo de confianza que se construye contiene el valor del parámetro especificado por H_0 , entonces ese valor es uno de los posibles valores del parámetro y H_0 no debe ser rechazada.
- Si el valor hipotético se encuentra fuera de los límites de confianza, la hipótesis nula es rechazada al nivel de significancia α .

1. Del libro Mendenhall resolver los ejercicios 9.18, 9.19 y 9.20(Mendenhall).
2. Del libro Mendenhall resolver los ejercicios: 9.23, 9.26 y 9.28.

47.22.3 2.2.3 Prueba de Hipótesis para una Proporción Binomial

Para una muestra aleatoria de n intentos idénticos, de una población binomial, la proporción muestral \hat{p} tiene una distribución aproximadamente normal cuando n es grande, con media p y error estándar

$$SE = \sqrt{\frac{pq}{n}}.$$

La prueba de hipótesis de la forma

$$\begin{aligned} H_0 &: p = p_0 \\ H_1 &: p > p_0, \text{ o } p < p_0 \text{ o } p \neq p_0 \end{aligned}$$

El estadístico de prueba se construye con el mejor estimador de la proporción verdadera, \hat{p} , con el estadístico de prueba z , que se distribuye normal estándar.

El procedimiento es

- 1) Hipótesis nula: $H_0 : p = p_0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis alternativa	$H_1 : p > p_0$ $H_1 : p < p_0$	$p \neq p_0$

- 3) Estadístico de prueba:

$$z = \frac{\hat{p} - p_0}{\sqrt{\frac{pq}{n}}}, \hat{p} = \frac{x}{n}$$

donde x es el número de éxitos en n intentos binomiales.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_0$ $z < -z_\alpha$ cuando $H_1 : p < p_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Ejemplo 40. A cualquier edad, alrededor del 20% de los adultos de cierto país realiza actividades de acondicionamiento físico al menos dos veces por semana. En una encuesta local de $n = 100$ adultos de más de 40 a nos, un total de 15 personas indicaron que realizaron actividad física al menos dos veces por semana. Estos datos indican que el porcentaje de participación para adultos de más de 40 a nos de edad es considerablemente menor a la cifra del 20%? Calcule el valor de p y úselo para sacar las conclusiones apropiadas.

1. Resolver los ejercicios: 9.30, 9.32, 9.33, 9.35 y 9.39.

47.22.4 2.2.4 Prueba de Hipótesis diferencia entre dos Proporciones Binomiales

Nota 81. Cuando se tienen dos muestras aleatorias independientes de dos poblaciones binomiales, el objetivo del experimento puede ser la diferencia $(p_1 - p_2)$ en las proporciones de individuos u objetos que poseen una característica específica en las dos poblaciones. En este caso se pueden utilizar los estimadores de las dos proporciones $(\hat{p}_1 - \hat{p}_2)$ con error estándar dado por

$$SE = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

considerando el estadístico z con un nivel de significancia $(1 - \alpha) 100\%$

Nota 82. La hipótesis nula a probarse es de la forma

$H_0: p_1 = p_2$ o equivalentemente $(p_1 - p_2) = 0$, contra una hipótesis alternativa H_1 de una o dos colas.

Nota 83. Para estimar el error estándar del estadístico z , se debe de utilizar el hecho de que suponiendo que H_0 es verdadera, las dos proporciones son iguales a algún valor común, p . Para obtener el mejor estimador de p es

$$p = \frac{\text{número total de éxitos}}{\text{Número total de pruebas}} = \frac{x_1 + x_2}{n_1 + n_2}$$

1) **Hipótesis Nula:** $H_0 : (p_1 - p_2) = 0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 :$	$H_1 : (p_1 - p_2) > 0$ $H_1 : (p_1 - p_2) < 0$	$H_1 : (p_1 - p_2) \neq 0$

3) Estadístico de prueba:

$$z = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} = \frac{(\hat{p}_1 - \hat{p}_2)}{\sqrt{\frac{pq}{n_1} + \frac{pq}{n_2}}}$$

donde $\hat{p}_1 = x_1/n_1$ y $\hat{p}_2 = x_2/n_2$, dado que el valor común para p_1 y p_2 es p , entonces $\hat{p} = \frac{x_1+x_2}{n_1+n_2}$ y por tanto el estadístico de prueba es

$$z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}\hat{q}\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_\alpha$ $z < -z_\alpha$ cuando $H_1 : p < p_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Ejemplo 41. Los registros de un hospital, indican que 52 hombres de una muestra de 1000 contra 23 mujeres de una muestra de 1000 fueron ingresados por enfermedad del corazón. Estos datos presentan suficiente evidencia para indicar un porcentaje más alto de enfermedades del corazón entre hombres ingresados al hospital?, utilizar distintos niveles de confianza de α .

1. Resolver los ejercicios 9.42
2. Resolver los ejercicios: 9.45, 9.48, 9.50

47.23 2.3 Muestras Pequeñas

47.23.1 2.3.1 Una media poblacional

1) **Hipótesis Nula:** $H_0 : \mu = \mu_0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 :$	$H_1 : \mu > \mu_0$ $H_1 : \mu < \mu_0$	$H_1 : \mu \neq \mu_0$

3) Estadístico de prueba:

$$t = \frac{\bar{x} - \mu_0}{\sqrt{\frac{s^2}{n}}}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Ejemplo 42. Las etiquetas en latas de un galón de pintura por lo general indican el tiempo de secado y el área puede cubrir una capa. Casi todas las marcas de pintura indican que, en una capa, un galón cubrirá entre 250 y 500 pies cuadrados, dependiendo de la textura de la superficie a pintarse, un fabricante, sin embargo afirma que un galón de su pintura cubrirá 400 pies cuadrados de área superficial. Para probar su afirmación, una muestra aleatoria de 10 latas de un galón de pintura blanca se empleó para pintar 10 áreas idénticas usando la misma clase de equipo. Las áreas reales en pies cuadrados cubiertas por estos 10 galones de pintura se dan a continuación:

310	311	412	368	447
376	303	410	365	350

Ejemplo 43. Los datos presentan suficiente evidencia para indicar que el promedio de la cobertura difiere de 400 pies cuadrados? encuentre el valor de p para la prueba y úselo para evaluar la significancia de los resultados.

1. Resolver los ejercicios: 10.2, 10.3, 10.5, 10.7, 10.9, 10.13 y 10.16

47.23.2 2.3.2 Diferencia entre dos medias poblacionales: M.A.I.

Nota 84. Cuando los tamaños de muestra son pequeños, no se puede asegurar que las medias muestrales sean normales, pero si las poblaciones originales son normales, entonces la distribución muestral de la diferencia de las medias muestrales, $(\bar{x}_1 - \bar{x}_2)$, será normal con media $(\mu_1 - \mu_2)$ y error estándar

$$ES = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

- 1) **Hipótesis Nula** $H_0 : (\mu_1 - \mu_2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\mu_1 - \mu_2) > D_0$ $H_1 : (\mu_1 - \mu_2) < D_0$	$H_1 : (\mu_1 - \mu_2) \neq D_0$

- 3) Estadístico de prueba:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - D_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

donde

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

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$z > z_0$ $z < -z_\alpha$ cuando $H_1 : (\mu_1 - \mu_2) < D_0$ cuando $p < \alpha$	$z > z_{\alpha/2}$ o $z < -z_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

47.23.3 2.3.3 Diferencia entre dos medias poblacionales: Diferencias Pareadas

- 1) **Hipótesis Nula:** $H_0 : \mu_d = 0$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: $H_1 : \mu_d$	$H_1 : \mu_d > 0$ $H_1 : \mu_d < 0$	$H_1 : \mu_d \neq 0$

- 3) Estadístico de prueba:

$$t = \frac{\bar{d}}{\sqrt{\frac{s_d^2}{n}}}$$

donde n es el número de diferencias pareadas, \bar{d} es la media de las diferencias muestrales, y s_d es la desviación estándar de las diferencias muestrales.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$t > t_\alpha$ $t < -t_\alpha$ cuando $H_1 : \mu < \mu_0$ cuando $p < \alpha$	$t > t_{\alpha/2}$ o $t < -t_{\alpha/2}$

Los valores críticos de t , $t_{-\alpha}$ y $t_{\alpha/2}$ están basados en $(n_1 + n_2 - 2)$ grados de libertad.

47.23.4 2.3.4 Inferencias con respecto a la Varianza Poblacional

- 1) **Hipótesis Nula:** $H_0 : \sigma^2 = \sigma_0^2$

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa: H_1	$H_1 : \sigma^2 > \sigma_0^2$ $H_1 : \sigma^2 < \sigma_0^2$	$H_1 : \sigma^2 \neq \sigma_0^2$

- 3) Estadístico de prueba:

$$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$$

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$\chi^2 > \chi_\alpha^2$ $\chi^2 < \chi_{(1-\alpha)}^2$ cuando $H_1 : \chi^2 < \chi_0^2$ cuando $p < \alpha$	$\chi^2 > \chi_{\alpha/2}^2$ o $\chi^2 < \chi_{(1-\alpha/2)}^2$

Los valores críticos de χ^2 , están basados en $(n_1 +)$ grados de libertad.

47.23.5 2.3.5 Comparación de dos varianzas poblacionales

- 1) **Hipótesis Nula** $H_0 : (\sigma_1^2 - \sigma_2^2) = D_0$,

donde D_0 es el valor, la diferencia, específico que se desea probar. En algunos casos se querrá demostrar que no hay diferencia alguna, es decir $D_0 = 0$.

	Prueba de una Cola	Prueba de dos colas
2) Hipótesis Alternativa	$H_1 : (\sigma_1^2 - \sigma_2^2) > D_0$ $H_1 : (\sigma_1^2 - \sigma_2^2) < D_0$	$H_1 : (\sigma_1^2 - \sigma_2^2) \neq D_0$

- 3) Estadístico de prueba:

$$F = \frac{s_1^2}{s_2^2}$$

donde s_1^2 es la varianza muestral más grande.

	Prueba de una Cola	Prueba de dos colas
4) Región de rechazo: rechazar H_0 cuando	$F > F_\alpha$ cuando $p < \alpha$	$F > F_{\alpha/2}$

47.24 Ejercicios

- Del libro Probabilidad y Estadística para Ingeniería de Hines, Montgomery, Goldsman y Borror resolver los siguientes ejercicios: 10-9, 10-10, 10-13, 10-16 y 10-20.
- Realizar un programa en R para cada una de las secciones y subsecciones revisadas en clase, para determinar intervalos de confianza.
- Aplicar los programas elaborados en el ejercicio anterior a la siguiente lista: 10-39, 10-41, 10-45, 10-47, 10-48, 10-50, 10-52, 10-54, 10-56, 10-57, 10-58, 10-65, 10-68, 10-72 y 10-73.
- Elaborar una rutina en R que grafique las siguientes distribuciones, permitiendo variar los parámetros de las distribuciones: Binomial, Uniforme continua, Gamma, Beta, Exponencial, Normal y t -Student.
- Presentar el primer capítulo del libro del curso en formato *Rnw* con su respectivo archivo *pdf* generado

47.25 Análisis de Regresión Lineal (RL)

Nota 85. — En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.
- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

47.26 Análisis de Regresión Lineal (RL)

Nota 86. – En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.
- A la relación entre estas variables se le denomina *modelo regresión de y en x_1, x_2, \dots, x_n* , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

47.26.1 Regresión Lineal Simple (RLS)

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (47.62)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (47.63)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación (??) se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (47.64)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (47.65)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (47.66)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (47.67)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (47.68)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (47.69)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (47.70)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (47.71)$$

47.26.2 Regresión Lineal Simple (RLS)

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (47.72)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (47.73)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación (??) se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (47.74)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (47.75)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (47.76)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (47.77)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (47.78)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (47.79)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (47.80)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (47.81)$$

47.27 3. Análisis de Regresión Lineal (RL)

Nota 87. – En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.

47.27.1 3.1 Regresión Lineal Simple (RLS)

- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (47.82)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

47.27.2 3.2 Método de Mínimos Cuadrados

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (47.83)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación ?? se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (47.84)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (47.85)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (47.86)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (47.87)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (47.88)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (47.89)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (47.90)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (47.91)$$

47.27.3 3.3 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Nota 88. – Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo.

- Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$\begin{aligned}
 E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right) \\
 &= \frac{1}{S_{xx}} E\left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x})\right) \\
 &= \frac{1}{S_{xx}} \left[\beta_0 E\left(\sum_{k=1}^n (x_k - \bar{x})\right) + E\left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x})\right) \right. \\
 &\quad \left. + E\left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x})\right) \right] = \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1
 \end{aligned}$$

por lo tanto

$$E(\hat{\beta}_1) = \beta_1 \quad (47.92)$$

Nota 89. Es decir, $\hat{\beta}_1$ es un estimador insesgado.

Ahora calculemos la varianza:

$$\begin{aligned}
 V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) \\
 &= \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) = \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 \\
 &= \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}}
 \end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad (47.93)$$

Proposición 15.

$$\begin{aligned}
 E(\hat{\beta}_0) &= \beta_0, \\
 V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \\
 Cov(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{\sigma^2 \bar{x}}{S_{xx}}.
 \end{aligned}$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$$e_k = y_k - \hat{y}_k, \text{ se le denomina } \mathbf{residuo}.$$

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \quad (47.94)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$\begin{aligned} SC_E &= \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \\ E(SC_E) &= (n-2)\sigma^2, \text{ por lo tanto} \\ \hat{\sigma}^2 &= \frac{SC_E}{n-2} = MC_E \text{ es un estimador insesgado de } \sigma^2. \end{aligned}$$

47.27.4 3.4 Prueba de Hipótesis en RLS

- Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza.
- Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$.

Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_1, \sigma^2)$.

De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??). Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (47.95)$$

que se distribuye t con $n-2$ grados de libertad bajo $H_0: \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (47.96)$$

Para β_0 se puede proceder de manera análoga para

$$H_0: \beta_0 = \beta_{0,0},$$

$$H_1: \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}, \quad (47.97)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (47.98)$$

- No rechazar $H_0 : \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0 : \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0 : \beta_1 = 0$ puede realizarse de la siguiente manera:

$$\begin{aligned} S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 \\ &= \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\ &= \sum_{k=1}^n \left[(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2 \right] \\ &= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ &\quad \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) \\ &= \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\ &= \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &= \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &= \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0. \end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (47.99)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (47.100)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (47.101)$$

Por lo tanto la ecuación (??) se puede reescribir como

$$S_{yy} = SC_R + SC_E \quad (47.102)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$

$$\begin{aligned} S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}) \\ S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R \end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 16.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx} \quad (47.103)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E}$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$.

El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n - 2$	MC_E	
Total	S_{yy}	$n - 1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}} \quad (47.104)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

47.27.5 Estimación de Intervalos en RLS

- Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros.
- El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión.

Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n - 2$ grados de libertad. Por tanto un intervalo de confianza de $100(1 - \alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}}. \quad (47.105)$$

De igual manera, para β_0 un intervalo de confianza al $100(1 - \alpha)\%$ es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (47.106)$$

47.28 3. Análisis de Regresión Lineal (RL)

Nota 90. – En muchos problemas hay dos o más variables relacionadas, para medir el grado de relación se utiliza el **análisis de regresión**.

- Supongamos que se tiene una única variable dependiente, y , y varias variables independientes, x_1, x_2, \dots, x_n .
- La variable y es una variable aleatoria, y las variables independientes pueden ser distribuidas independiente o conjuntamente.

47.28.1 3.1 Regresión Lineal Simple (RLS)

- A la relación entre estas variables se le denomina modelo regresión de y en x_1, x_2, \dots, x_n , por ejemplo $y = \phi(x_1, x_2, \dots, x_n)$, lo que se busca es una función que mejor aproxime a $\phi(\cdot)$.

Supongamos que de momento solamente se tienen una variable independiente x , para la variable de respuesta y . Y supongamos que la relación que hay entre x y y es una línea recta, y que para cada observación de x , y es una variable aleatoria.

El valor esperado de y para cada valor de x es

$$E(y|x) = \beta_0 + \beta_1 x \quad (47.107)$$

β_0 es la ordenada al origen y β_1 la pendiente de la recta en cuestión, ambas constantes desconocidas.

47.28.2 3.2 Método de Mínimos Cuadrados

Supongamos que cada observación y se puede describir por el modelo

$$y = \beta_0 + \beta_1 x + \epsilon \quad (47.108)$$

donde ϵ es un error aleatorio con media cero y varianza σ^2 . Para cada valor y_i se tiene ϵ_i variables aleatorias no correlacionadas, cuando se incluyen en el modelo ??, este se le llama *modelo de regresión lineal simple*.

Suponga que se tienen n pares de observaciones $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, estos datos pueden utilizarse para estimar los valores de β_0 y β_1 . Esta estimación es por el **métodos de mínimos cuadrados**.

Entonces la ecuación ?? se puede reescribir como

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \text{ para } i = 1, 2, \dots, n. \quad (47.109)$$

Si consideramos la suma de los cuadrados de los errores aleatorios, es decir, el cuadrado de la diferencia entre las observaciones con la recta de regresión

$$L = \sum_{i=1}^n \epsilon^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \quad (47.110)$$

Para obtener los estimadores por mínimos cuadrados de β_0 y β_1 , $\hat{\beta}_0$ y $\hat{\beta}_1$, es preciso calcular las derivadas parciales con respecto a β_0 y β_1 , igualar a cero y resolver el sistema de ecuaciones lineales resultante:

$$\begin{aligned} \frac{\partial L}{\partial \beta_0} &= 0 \\ \frac{\partial L}{\partial \beta_1} &= 0 \end{aligned}$$

evaluando en $\hat{\beta}_0$ y $\hat{\beta}_1$, se tiene

$$\begin{aligned} -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) &= 0 \\ -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i &= 0 \end{aligned}$$

simplificando

$$\begin{aligned} n\hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i &= \sum_{i=1}^n y_i \\ \hat{\beta}_0 \sum_{i=1}^n x_i + \hat{\beta}_1 \sum_{i=1}^n x_i^2 &= \sum_{i=1}^n x_i y_i \end{aligned}$$

Las ecuaciones anteriores se les denominan *ecuaciones normales de mínimos cuadrados* con solución

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (47.111)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_{i=1}^n y_i) (\sum_{i=1}^n x_i)}{\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2} \quad (47.112)$$

entonces el modelo de regresión lineal simple ajustado es

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad (47.113)$$

Se introduce la siguiente notación

$$S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \quad (47.114)$$

$$S_{xy} = \sum_{i=1}^n y_i (x_i - \bar{x}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \quad (47.115)$$

y por tanto

$$\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} \quad (47.116)$$

47.28.3 3.3 Propiedades de los Estimadores $\hat{\beta}_0$ y $\hat{\beta}_1$

Nota 91. – Las propiedades estadísticas de los estimadores de mínimos cuadrados son útiles para evaluar la suficiencia del modelo.

– Dado que $\hat{\beta}_0$ y $\hat{\beta}_1$ son combinaciones lineales de las variables aleatorias y_i , también resultan ser variables aleatorias.

A saber

$$\begin{aligned} E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}} E\left(\sum_{i=1}^n y_i (x_i - \bar{x})\right) \\ &= \frac{1}{S_{xx}} E\left(\sum_{i=1}^n (\beta_0 + \beta_1 x_i + \epsilon_i) (x_i - \bar{x})\right) \\ &= \frac{1}{S_{xx}} \left[\beta_0 E\left(\sum_{k=1}^n (x_k - \bar{x})\right) + E\left(\beta_1 \sum_{k=1}^n x_k (x_k - \bar{x})\right) \right. \\ &\quad \left. + E\left(\sum_{k=1}^n \epsilon_k (x_k - \bar{x})\right) \right] = \frac{1}{S_{xx}} \beta_1 S_{xx} = \beta_1 \end{aligned}$$

por lo tanto

$$E(\hat{\beta}_1) = \beta_1 \quad (47.117)$$

Nota 92. Es decir, $\hat{\beta}_1$ es un estimador insesgado.

Ahora calculemos la varianza:

$$\begin{aligned} V(\hat{\beta}_1) &= V\left(\frac{S_{xy}}{S_{xx}}\right) = \frac{1}{S_{xx}^2} V\left(\sum_{k=1}^n y_k (x_k - \bar{x})\right) \\ &= \frac{1}{S_{xx}^2} \sum_{k=1}^n V(y_k (x_k - \bar{x})) = \frac{1}{S_{xx}^2} \sum_{k=1}^n \sigma^2 (x_k - \bar{x})^2 \\ &= \frac{\sigma^2}{S_{xx}^2} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{\sigma^2}{S_{xx}} \end{aligned}$$

por lo tanto

$$V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}} \quad (47.118)$$

Proposición 17.

$$\begin{aligned}E(\hat{\beta}_0) &= \beta_0, \\V(\hat{\beta}_0) &= \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right), \\Cov(\hat{\beta}_0, \hat{\beta}_1) &= -\frac{\sigma^2 \bar{x}}{S_{xx}}.\end{aligned}$$

Para estimar σ^2 es preciso definir la diferencia entre la observación y_k , y el valor predicho \hat{y}_k , es decir

$$e_k = y_k - \hat{y}_k, \text{ se le denomina } \mathbf{residuo}.$$

La suma de los cuadrados de los errores de los residuos, *suma de cuadrados del error*

$$SC_E = \sum_{k=1}^n e_k^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \quad (47.119)$$

sustituyendo $\hat{y}_k = \hat{\beta}_0 + \hat{\beta}_1 x_k$ se obtiene

$$\begin{aligned}SC_E &= \sum_{k=1}^n y_k^2 - n\bar{y}^2 - \hat{\beta}_1 S_{xy} = S_{yy} - \hat{\beta}_1 S_{xy}, \\E(SC_E) &= (n-2)\sigma^2, \text{ por lo tanto} \\ \hat{\sigma}^2 &= \frac{SC_E}{n-2} = MC_E \text{ es un estimador insesgado de } \sigma^2.\end{aligned}$$

47.28.4 3.4 Prueba de Hipótesis en RLS

- Para evaluar la suficiencia del modelo de regresión lineal simple, es necesario llevar a cabo una prueba de hipótesis respecto de los parámetros del modelo así como de la construcción de intervalos de confianza.
- Para poder realizar la prueba de hipótesis sobre la pendiente y la ordenada al origen de la recta de regresión es necesario hacer el supuesto de que el error ϵ_i se distribuye normalmente, es decir $\epsilon_i \sim N(0, \sigma^2)$.

Suponga que se desea probar la hipótesis de que la pendiente es igual a una constante, $\beta_{0,1}$ las hipótesis Nula y Alternativa son:

$$H_0: \beta_1 = \beta_{1,0},$$

$$H_1: \beta_1 \neq \beta_{1,0}.$$

donde dado que las $\epsilon_i \sim N(0, \sigma^2)$, se tiene que y_i son variables aleatorias normales $N(\beta_0 + \beta_1 x_1, \sigma^2)$.

De las ecuaciones (??) se desprende que $\hat{\beta}_1$ es combinación lineal de variables aleatorias normales independientes, es decir, $\hat{\beta}_1 \sim N(\beta_1, \sigma^2/S_{xx})$, recordar las ecuaciones (??) y (??).

Entonces se tiene que el estadístico de prueba apropiado es

$$t_0 = \frac{\hat{\beta}_1 - \hat{\beta}_{1,0}}{\sqrt{MC_E/S_{xx}}} \quad (47.120)$$

que se distribuye t con $n-2$ grados de libertad bajo $H_0: \beta_1 = \beta_{1,0}$. Se rechaza H_0 si

$$|t_0| > t_{\alpha/2, n-2}. \quad (47.121)$$

Para β_0 se puede proceder de manera análoga para

$$H_0 : \beta_0 = \beta_{0,0},$$

$$H_1 : \beta_0 \neq \beta_{0,0},$$

con $\hat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)\right)$, por lo tanto

$$t_0 = \frac{\hat{\beta}_0 - \beta_{0,0}}{MC_E \left[\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right]}, \quad (47.122)$$

con el que rechazamos la hipótesis nula si

$$|t_0| > t_{\alpha/2, n-2}. \quad (47.123)$$

- No rechazar $H_0 : \beta_1 = 0$ es equivalente a decir que no hay relación lineal entre x y y .
- Alternativamente, si $H_0 : \beta_1 = 0$ se rechaza, esto implica que x explica la variabilidad de y , es decir, podría significar que la línea recta es el modelo adecuado.

El procedimiento de prueba para $H_0 : \beta_1 = 0$ puede realizarse de la siguiente manera:

$$\begin{aligned} S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ S_{yy} &= \sum_{k=1}^n (y_k - \bar{y})^2 = \sum_{k=1}^n (y_k - \hat{y}_k + \hat{y}_k - \bar{y})^2 \\ &= \sum_{k=1}^n [(\hat{y}_k - \bar{y}) + (y_k - \hat{y}_k)]^2 \\ &= \sum_{k=1}^n \left[(\hat{y}_k - \bar{y})^2 + 2(\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + (y_k - \hat{y}_k)^2 \right] \\ &= \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + 2 \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) + \sum_{k=1}^n (y_k - \hat{y}_k)^2 \\ &= \sum_{k=1}^n (\hat{y}_k - \bar{y})(y_k - \hat{y}_k) = \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \sum_{k=1}^n \bar{y} (y_k - \hat{y}_k) \\ &= \sum_{k=1}^n \hat{y}_k (y_k - \hat{y}_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{y}_k) \\ &= \sum_{k=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_k) (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &= \sum_{k=1}^n \hat{\beta}_0 (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \sum_{k=1}^n \hat{\beta}_1 x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &= \hat{\beta}_0 \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) + \hat{\beta}_1 \sum_{k=1}^n x_k (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) \\ &\quad - \bar{y} \sum_{k=1}^n (y_k - \hat{\beta}_0 - \hat{\beta}_1 x_k) = 0 + 0 + 0 = 0. \end{aligned}$$

Por lo tanto, efectivamente se tiene

$$S_{yy} = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 + \sum_{k=1}^n (y_k - \hat{y}_k)^2, \quad (47.124)$$

donde se hacen las definiciones

$$SC_E = \sum_{k=1}^n (\hat{y}_k - \bar{y})^2 \cdots \text{Suma de Cuadrados del Error} \quad (47.125)$$

$$SC_R = \sum_{k=1}^n (y_k - \hat{y}_k)^2 \cdots \text{Suma de Regresión de Cuadrados} \quad (47.126)$$

Por lo tanto la ecuación (??) se puede reescribir como

$$S_{yy} = SC_R + SC_E \quad (47.127)$$

recordemos que $SC_E = S_{yy} - \hat{\beta}_1 S_{xy}$

$$\begin{aligned} S_{yy} &= SC_R + (S_{yy} - \hat{\beta}_1 S_{xy}) \\ S_{xy} &= \frac{1}{\hat{\beta}_1} SC_R \end{aligned}$$

S_{xy} tiene $n - 1$ grados de libertad y SC_R y SC_E tienen 1 y $n - 2$ grados de libertad respectivamente.

Proposición 18.

$$E(SC_R) = \sigma^2 + \beta_1 S_{xx} \quad (47.128)$$

además, SC_E y SC_R son independientes.

Recordemos que $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$. Para $H_0 : \beta_1 = 0$ verdadera,

$$F_0 = \frac{SC_R/1}{SC_E/(n-2)} = \frac{MC_R}{MC_E}$$

se distribuye $F_{1,n-2}$, y se rechazaría H_0 si $F_0 > F_{\alpha,1,n-2}$.

El procedimiento de prueba de hipótesis puede presentarse como la tabla de análisis de varianza siguiente

Fuente de variación	Suma de Cuadrados	Grados de Libertad	Media Cuadrática	F_0
Regresión	SC_R	1	MC_R	MC_R/MC_E
Error Residual	SC_E	$n - 2$	MC_E	
Total	S_{yy}	$n - 1$		

La prueba para la significación de la regresión puede desarrollarse basándose en la expresión (??), con $\hat{\beta}_{1,0} = 0$, es decir

$$t_0 = \frac{\hat{\beta}_1}{\sqrt{MC_E/S_{xx}}} \quad (47.129)$$

Elevando al cuadrado ambos términos:

$$t_0^2 = \frac{\hat{\beta}_1^2 S_{xx}}{MC_E} = \frac{\hat{\beta}_1 S_{xy}}{MC_E} = \frac{MC_R}{MC_E}$$

Observar que $t_0^2 = F_0$, por tanto la prueba que se utiliza para t_0 es la misma que para F_0 .

47.28.5 Estimación de Intervalos en RLS

- Además de la estimación puntual para los parámetros β_1 y β_0 , es posible obtener estimaciones del intervalo de confianza de estos parámetros.
- El ancho de estos intervalos de confianza es una medida de la calidad total de la recta de regresión.

Si los ϵ_k se distribuyen normal e independientemente, entonces

$$\frac{(\hat{\beta}_1 - \beta_1)}{\sqrt{\frac{MC_E}{S_{xx}}}} \quad y \quad \frac{(\hat{\beta}_0 - \beta_0)}{\sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}}$$

se distribuyen t con $n - 2$ grados de libertad. Por tanto un intervalo de confianza de $100(1 - \alpha)\%$ para β_1 está dado por

$$\hat{\beta}_1 - t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}} \leq \beta_1 \leq \hat{\beta}_1 + t_{\alpha/2, n-2} \sqrt{\frac{MC_E}{S_{xx}}}. \quad (47.130)$$

De igual manera, para β_0 un intervalo de confianza al $100(1 - \alpha)\%$ es

$$\hat{\beta}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \leq \beta_0 \leq \hat{\beta}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \quad (47.131)$$

47.28.6 Predicción

Supongamos que se tiene un valor x_0 de interés, entonces la estimación puntual de este nuevo valor

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 \quad (47.132)$$

Nota 93. Esta nueva observación es independiente de las utilizadas para obtener el modelo de regresión, por tanto, el intervalo en torno a la recta de regresión es inapropiado, puesto que se basa únicamente en los datos empleados para ajustar el modelo de regresión.

El intervalo de confianza en torno a la recta de regresión se refiere a la respuesta media verdadera $x = x_0$, no a observaciones futuras.

Sea y_0 la observación futura en $x = x_0$, y sea \hat{y}_0 dada en la ecuación anterior, el estimador de y_0 . Si se define la variable aleatoria

$$w = y_0 - \hat{y}_0,$$

esta se distribuye normalmente con media cero y varianza

$$V(w) = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]$$

dado que y_0 es independiente de \hat{y}_0 , por lo tanto el intervalo de predicción al nivel α para futuras observaciones x_0 es

$$\begin{aligned} \hat{y}_0 - t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]} &\leq y_0 \\ &\leq \hat{y}_0 + t_{\alpha/2, n-2} \sqrt{MC_E \left[1 + \frac{1}{n} + \frac{(x - x_0)^2}{S_{xx}} \right]}. \end{aligned}$$

47.28.7 Prueba de falta de ajuste

Es común encontrar que el modelo ajustado no satisface totalmente el modelo necesario para los datos, en este caso es preciso saber qué tan bueno es el modelo propuesto. Para esto se propone la siguiente prueba de hipótesis:

H_0 : El modelo propuesto se ajusta adecuadamente a los datos.

H_1 : El modelo NO se ajusta a los datos.

La prueba implica dividir la suma de cuadrados del error o del residuo en las siguientes dos componentes:

$$SC_E = SC_{EP} + SC_{FDA}$$

donde SC_{EP} es la suma de cuadrados atribuibles al error puro, y SC_{FDA} es la suma de cuadrados atribuible a la falta de ajuste del modelo.

47.28.8 Coeficiente de Determinación

La cantidad

$$R^2 = \frac{SC_R}{S_{yy}} = 1 - \frac{SC_E}{S_{yy}} \quad (47.133)$$

se denomina coeficiente de determinación y se utiliza para saber si el modelo de regresión es suficiente o no. Se puede demostrar que $0 \leq R^2 \leq 1$, una manera de interpretar este valor es que si $R^2 = k$, entonces el modelo de regresión explica el $k * 100\%$ de la variabilidad en los datos. R^2

- No mide la magnitud de la pendiente de la recta de regresión
- Un valor grande de R^2 no implica una pendiente empinada.
- No mide la suficiencia del modelo.
- Valores grandes de R^2 no implican necesariamente que el modelo de regresión proporcionará predicciones precisas para futuras observaciones.

47.29 Introducción

La preparación de datos y la selección de variables son pasos cruciales en el proceso de modelado estadístico. Un modelo bien preparado y con las variables adecuadas puede mejorar significativamente la precisión y la interpretabilidad del modelo. Este capítulo proporciona una revisión detallada de las técnicas de limpieza de datos, tratamiento de datos faltantes, codificación de variables categóricas y selección de variables.

47.30 Importancia de la Preparación de Datos

La calidad de los datos es fundamental para el éxito de cualquier análisis estadístico. Los datos sin limpiar pueden llevar a modelos inexactos y conclusiones erróneas. La preparación de datos incluye varias etapas:

- Limpieza de datos
- Tratamiento de datos faltantes
- Codificación de variables categóricas
- Selección y transformación de variables

47.31 Limpieza de Datos

La limpieza de datos es el proceso de detectar y corregir (o eliminar) los datos incorrectos, incompletos o irrelevantes. Este proceso incluye:

- Eliminación de duplicados
- Corrección de errores tipográficos
- Consistencia de formato
- Tratamiento de valores extremos (outliers)

47.32 Tratamiento de Datos Faltantes

Los datos faltantes son un problema común en los conjuntos de datos y pueden afectar la calidad de los modelos. Hay varias estrategias para manejar los datos faltantes:

- **Eliminación de Datos Faltantes:** Se eliminan las filas o columnas con datos faltantes.
- **Imputación:** Se reemplazan los valores faltantes con estimaciones, como la media, la mediana o la moda.
- **Modelos Predictivos:** Se utilizan modelos predictivos para estimar los valores faltantes.

47.32.1 Imputación de la Media

Una técnica común es reemplazar los valores faltantes con la media de la variable. Esto se puede hacer de la siguiente manera:

$$x_i = \begin{cases} x_i & \text{si } x_i \text{ no es faltante} \\ \bar{x} & \text{si } x_i \text{ es faltante} \end{cases}$$

donde \bar{x} es la media de la variable.

47.33 Codificación de Variables Categóricas

Las variables categóricas deben ser convertidas a un formato numérico antes de ser usadas en un modelo de regresión logística. Hay varias técnicas para codificar variables categóricas:

47.33.1 Codificación One-Hot

La codificación one-hot crea una columna binaria para cada categoría. Por ejemplo, si tenemos una variable categórica con tres categorías (A, B, C), se crean tres columnas:

$$\begin{aligned} A &= [1, 0, 0] \\ B &= [0, 1, 0] \\ C &= [0, 0, 1] \end{aligned}$$

47.33.2 Codificación Ordinal

La codificación ordinal asigna un valor entero único a cada categoría, preservando el orden natural de las categorías. Por ejemplo:

Bajo	=	1
Medio	=	2
Alto	=	3

47.34 Selección de Variables

La selección de variables es el proceso de elegir las variables más relevantes para el modelo. Existen varias técnicas para la selección de variables:

47.34.1 Métodos de Filtrado

Los métodos de filtrado seleccionan variables basadas en criterios estadísticos, como la correlación o la chi-cuadrado. Algunas técnicas comunes incluyen:

- **Análisis de Correlación:** Se seleccionan variables con alta correlación con la variable dependiente y baja correlación entre ellas.
- **Pruebas de Chi-cuadrado:** Se utilizan para variables categóricas para determinar la asociación entre la variable independiente y la variable dependiente.

47.34.2 Métodos de Wrapper

Los métodos de wrapper evalúan múltiples combinaciones de variables y seleccionan la combinación que optimiza el rendimiento del modelo. Ejemplos incluyen:

- **Selección hacia Adelante:** Comienza con un modelo vacío y agrega variables una por una, seleccionando la variable que mejora más el modelo en cada paso.
- **Selección hacia Atrás:** Comienza con todas las variables y elimina una por una, removiendo la variable que tiene el menor impacto en el modelo en cada paso.
- **Selección Paso a Paso:** Combina la selección hacia adelante y hacia atrás, agregando y eliminando variables según sea necesario.

47.34.3 Métodos Basados en Modelos

Los métodos basados en modelos utilizan técnicas de regularización como Lasso y Ridge para seleccionar variables. Estas técnicas añaden un término de penalización a la función de costo para evitar el sobreajuste.

Regresión Lasso

La regresión Lasso (Least Absolute Shrinkage and Selection Operator) añade una penalización L_1 a la función de costo:

$$J(\beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

donde λ es el parámetro de regularización que controla la cantidad de penalización.

Regresión Ridge

La regresión Ridge añade una penalización L_2 a la función de costo:

$$J(\beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

donde λ es el parámetro de regularización.

47.35 Implementación en R

47.35.1 Limpieza de Datos

Para ilustrar la limpieza de datos en R, considere el siguiente conjunto de datos:

```
data <- data.frame(
  var1 = c(1, 2, 3, NA, 5),
  var2 = c("A", "B", "A", "B", "A"),
  var3 = c(10, 15, 10, 20, 25)
)

# Eliminación de filas con datos faltantes
data_clean <- na.omit(data)

# Imputación de la media
data$var1[is.na(data$var1)] <- mean(data$var1, na.rm = TRUE)
```

47.35.2 Codificación de Variables Categóricas

Para codificar variables categóricas, utilice la función ‘model.matrix’:

```
data <- data.frame(
  var1 = c(1, 2, 3, 4, 5),
  var2 = c("A", "B", "A", "B", "A")
)

# Codificación one-hot
data_onehot <- model.matrix(~ var2 - 1, data = data)
```

47.35.3 Selección de Variables

Para la selección de variables, utilice el paquete ‘caret’:

```
library(caret)

# Dividir los datos en conjuntos de entrenamiento y prueba
set.seed(123)
trainIndex <- createDataPartition(data$var1, p = .8,
                                   list = FALSE,
                                   times = 1)

dataTrain <- data[trainIndex,]
dataTest <- data[-trainIndex,]

# Modelo de regresión logística
model <- train(var1 ~ ., data = dataTrain, method = "glm", family = "binomial")

# Selección de variables
model <- stepAIC(model, direction = "both")
summary(model)
```

47.36 Introducción

Evaluar la calidad y el rendimiento de un modelo de regresión logística es crucial para asegurar que las predicciones sean precisas y útiles. Este capítulo se centra en las técnicas y métricas utilizadas para evaluar modelos de clasificación binaria, así como en la validación cruzada, una técnica para evaluar la generalización del modelo.

47.37 Métricas de Evaluación del Modelo

Las métricas de evaluación permiten cuantificar la precisión y el rendimiento de un modelo. Algunas de las métricas más comunes incluyen:

47.37.1 Curva ROC y AUC

La curva ROC (Receiver Operating Characteristic) es una representación gráfica de la sensibilidad (verdaderos positivos) frente a 1 - especificidad (falsos positivos). El área bajo la curva (AUC) mide la capacidad del modelo para distinguir entre las clases.

$$\begin{aligned}\text{Sensibilidad} &= \frac{\text{TP}}{\text{TP} + \text{FN}} \\ \text{Especificidad} &= \frac{\text{TN}}{\text{TN} + \text{FP}}\end{aligned}$$

47.37.2 Matriz de Confusión

La matriz de confusión es una tabla que muestra el rendimiento del modelo comparando las predicciones con los valores reales. Los términos incluyen:

- **Verdaderos Positivos (TP)**: Predicciones correctas de la clase positiva.
- **Falsos Positivos (FP)**: Predicciones incorrectas de la clase positiva.
- **Verdaderos Negativos (TN)**: Predicciones correctas de la clase negativa.
- **Falsos Negativos (FN)**: Predicciones incorrectas de la clase negativa.

	Predicción Positiva	Predicción Negativa
Real Positiva	TP	FN
Real Negativa	FP	TN

Table 47.3: Matriz de Confusión

47.37.3 Precisión, Recall y F1-Score

$$\begin{aligned}\text{Precisión} &= \frac{\text{TP}}{\text{TP} + \text{FP}} \\ \text{Recall} &= \frac{\text{TP}}{\text{TP} + \text{FN}} \\ \text{F1-Score} &= 2 \cdot \frac{\text{Precisión} \cdot \text{Recall}}{\text{Precisión} + \text{Recall}}\end{aligned}$$

47.37.4 Log-Loss

La pérdida logarítmica (Log-Loss) mide la precisión de las probabilidades predichas. La fórmula es:

$$\text{Log-Loss} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

donde y_i son los valores reales y p_i son las probabilidades predichas.

47.38 Validación Cruzada

La validación cruzada es una técnica para evaluar la capacidad de generalización de un modelo. Existen varios tipos de validación cruzada:

47.38.1 K-Fold Cross-Validation

En K-Fold Cross-Validation, los datos se dividen en K subconjuntos. El modelo se entrena K veces, cada vez utilizando K-1 subconjuntos para el entrenamiento y el subconjunto restante para la validación.

$$\text{Error Medio} = \frac{1}{K} \sum_{k=1}^K \text{Error}_k$$

47.38.2 Leave-One-Out Cross-Validation (LOOCV)

En LOOCV, cada observación se usa una vez como conjunto de validación y las restantes como conjunto de entrenamiento. Este método es computacionalmente costoso pero útil para conjuntos de datos pequeños.

47.39 Ajuste y Sobreajuste del Modelo

El ajuste adecuado del modelo es crucial para evitar el sobreajuste (overfitting) y el subajuste (underfitting).

47.39.1 Sobreajuste

El sobreajuste ocurre cuando un modelo se ajusta demasiado bien a los datos de entrenamiento, capturando ruido y patrones irrelevantes. Los síntomas incluyen una alta precisión en el entrenamiento y baja precisión en la validación.

47.39.2 Subajuste

El subajuste ocurre cuando un modelo no captura los patrones subyacentes de los datos. Los síntomas incluyen baja precisión tanto en el entrenamiento como en la validación.

47.39.3 Regularización

La regularización es una técnica para prevenir el sobreajuste añadiendo un término de penalización a la función de costo. Las técnicas comunes incluyen:

- Regresión Lasso (L1)
- Regresión Ridge (L2)

47.40 Implementación en R

47.40.1 Evaluación del Modelo

```
# Cargar el paquete necesario
library(caret)

# Dividir los datos en conjuntos de entrenamiento y prueba
set.seed(123)
trainIndex <- createDataPartition(data$var1, p = .8,
                                   list = FALSE,
                                   times = 1)

dataTrain <- data[trainIndex,]
dataTest <- data[-trainIndex,]
```

```
# Entrenar el modelo de regresión logística
model <- train(var1 ~ ., data = dataTrain, method = "glm", family = "binomial")

# Predicciones en el conjunto de prueba
predictions <- predict(model, dataTest)

# Matriz de confusión
confusionMatrix(predictions, dataTest$var1)
```

47.40.2 Validación Cruzada

```
# K-Fold Cross-Validation
control <- trainControl(method = "cv", number = 10)
model_cv <- train(var1 ~ ., data = dataTrain, method = "glm",
                  family = "binomial", trControl = control)

# Evaluación del modelo con validación cruzada
print(model_cv)
```

47.41 Introducción

El diagnóstico del modelo y el ajuste de parámetros son pasos esenciales para mejorar la precisión y la robustez de los modelos de regresión logística. Este capítulo se enfoca en las técnicas para diagnosticar problemas en los modelos y en métodos para ajustar los parámetros de manera óptima.

47.42 Diagnóstico del Modelo

El diagnóstico del modelo implica evaluar el rendimiento del modelo y detectar posibles problemas, como el sobreajuste, la multicolinealidad y la influencia de puntos de datos individuales.

47.42.1 Residuos

Los residuos son las diferencias entre los valores observados y los valores predichos por el modelo. El análisis de residuos puede revelar patrones que indican problemas con el modelo.

$$\text{Residuo}_i = y_i - \hat{y}_i$$

Residuos Estudiantizados

Los residuos estudiantizados se ajustan por la variabilidad del residuo y se utilizan para detectar outliers.

$$r_i = \frac{\text{Residuo}_i}{\hat{\sigma}\sqrt{1 - h_i}}$$

donde h_i es el leverage del punto de datos.

47.42.2 Influencia

La influencia mide el impacto de un punto de datos en los coeficientes del modelo. Los puntos con alta influencia pueden distorsionar el modelo.

Distancia de Cook

La distancia de Cook es una medida de la influencia de un punto de datos en los coeficientes del modelo.

$$D_i = \frac{r_i^2}{p} \cdot \frac{h_i}{1 - h_i}$$

donde p es el número de parámetros en el modelo.

47.42.3 Multicolinealidad

La multicolinealidad ocurre cuando dos o más variables independientes están altamente correlacionadas. Esto puede inflar las varianzas de los coeficientes y hacer que el modelo sea inestable.

Factor de Inflación de la Varianza (VIF)

El VIF mide cuánto se inflan las varianzas de los coeficientes debido a la multicolinealidad.

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

donde R_j^2 es el coeficiente de determinación de la regresión de la variable j contra todas las demás variables.

47.43 Ajuste de Parámetros

El ajuste de parámetros implica seleccionar los valores óptimos para los hiperparámetros del modelo. Esto puede mejorar el rendimiento y prevenir el sobreajuste.

47.43.1 Grid Search

El grid search es un método exhaustivo para ajustar los parámetros. Se define una rejilla de posibles valores de parámetros y se evalúa el rendimiento del modelo para cada combinación.

47.43.2 Random Search

El random search selecciona aleatoriamente combinaciones de valores de parámetros dentro de un rango especificado. Es menos exhaustivo que el grid search, pero puede ser más eficiente.

47.43.3 Bayesian Optimization

La optimización bayesiana utiliza modelos probabilísticos para seleccionar iterativamente los valores de parámetros más prometedores.

47.44 Implementación en R

47.44.1 Diagnóstico del Modelo

```
# Cargar el paquete necesario
library(car)

# Residuos estudentizados
dataTrain$resid <- rstudent(model)
hist(dataTrain$resid, breaks = 20, main = "Residuos Estudentizados")

# Distancia de Cook
dataTrain$cook <- cooks.distance(model)
plot(dataTrain$cook, type = "h", main = "Distancia de Cook")

# Factor de Inflación de la Varianza
vif_values <- vif(model)
print(vif_values)
```

47.44.2 Ajuste de Parámetros

```
# Grid Search con caret
control <- trainControl(method = "cv", number = 10)
tune_grid <- expand.grid(.alpha = c(0, 0.5, 1), .lambda = seq(0.01, 0.1, by = 0.01))

model_tune <- train(var1 ~ ., data = dataTrain, method = "glmnet",
                    trControl = control, tuneGrid = tune_grid)

print(model_tune)
```

47.45 Introducción

Interpretar correctamente los resultados de un modelo de regresión logística es esencial para tomar decisiones informadas. Este capítulo se centra en la interpretación de los coeficientes del modelo, los odds ratios, los intervalos de confianza y la significancia estadística.

47.46 Coeficientes de Regresión Logística

Los coeficientes de regresión logística representan la relación entre las variables independientes y la variable dependiente en términos de log-odds.

47.46.1 Interpretación de los Coeficientes

Cada coeficiente β_j en el modelo de regresión logística se interpreta como el cambio en el log-odds de la variable dependiente por unidad de cambio en la variable independiente X_j .

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

47.46.2 Signo de los Coeficientes

- **Coeficiente Positivo:** Un coeficiente positivo indica que un aumento en la variable independiente está asociado con un aumento en el log-odds de la variable dependiente.
- **Coeficiente Negativo:** Un coeficiente negativo indica que un aumento en la variable independiente está asociado con una disminución en el log-odds de la variable dependiente.

47.47 Odds Ratios

Las odds ratios proporcionan una interpretación más intuitiva de los coeficientes de regresión logística. La odds ratio para una variable independiente X_j se calcula como e^{β_j} .

47.47.1 Cálculo de las Odds Ratios

$$OR_j = e^{\beta_j}$$

47.47.2 Interpretación de las Odds Ratios

- **OR > 1:** Un OR mayor que 1 indica que un aumento en la variable independiente está asociado con un aumento en las odds de la variable dependiente.
- **OR < 1:** Un OR menor que 1 indica que un aumento en la variable independiente está asociado con una disminución en las odds de la variable dependiente.
- **OR = 1:** Un OR igual a 1 indica que la variable independiente no tiene efecto sobre las odds de la variable dependiente.

47.48 Intervalos de Confianza

Los intervalos de confianza proporcionan una medida de la incertidumbre asociada con los estimadores de los coeficientes. Un intervalo de confianza del 95% para un coeficiente β_j indica que, en el 95% de las muestras, el intervalo contendrá el valor verdadero de β_j .

47.48.1 Cálculo de los Intervalos de Confianza

Para calcular un intervalo de confianza del 95% para un coeficiente β_j , utilizamos la fórmula:

$$\beta_j \pm 1.96 \cdot SE(\beta_j)$$

donde $SE(\beta_j)$ es el error estándar de β_j .

47.49 Significancia Estadística

La significancia estadística se utiliza para determinar si los coeficientes del modelo son significativamente diferentes de cero. Esto se evalúa mediante pruebas de hipótesis.

47.49.1 Prueba de Hipótesis

Para cada coeficiente β_j , la hipótesis nula H_0 es que $\beta_j = 0$. La hipótesis alternativa H_a es que $\beta_j \neq 0$.

47.49.2 P-valor

El p-valor indica la probabilidad de obtener un coeficiente tan extremo como el observado, asumiendo que la hipótesis nula es verdadera. Un p-valor menor que el nivel de significancia α (típicamente 0.05) indica que podemos rechazar la hipótesis nula.

47.50 Implementación en R

47.50.1 Cálculo de Coeficientes y Odds Ratios

```
# Cargar el paquete necesario
library(broom)

# Entrenar el modelo de regresión logística
model <- glm(var1 ~ ., data = dataTrain, family = "binomial")

# Coeficientes del modelo
coef(model)

# Odds ratios
exp(coef(model))
```

47.50.2 Intervalos de Confianza

```
# Intervalos de confianza para los coeficientes
confint(model)

# Intervalos de confianza para las odds ratios
exp(confint(model))
```

47.50.3 P-valores y Significancia Estadística

```
# Resumen del modelo con p-valores
summary(model)
```

47.51 Introducción

La regresión logística multinomial y el análisis de supervivencia son extensiones de la regresión logística binaria. Este capítulo se enfoca en las técnicas y aplicaciones de estos métodos avanzados.

47.52 Regresión Logística Multinomial

La regresión logística multinomial se utiliza cuando la variable dependiente tiene más de dos categorías.

47.52.1 Modelo Multinomial

El modelo de regresión logística multinomial generaliza el modelo binario para manejar múltiples categorías. La probabilidad de que una observación pertenezca a la categoría k se expresa como:

$$P(Y = k) = \frac{e^{\beta_{0k} + \beta_{1k}X_1 + \dots + \beta_{nk}X_n}}{\sum_{j=1}^K e^{\beta_{0j} + \beta_{1j}X_1 + \dots + \beta_{nj}X_n}}$$

47.52.2 Estimación de Parámetros

Los coeficientes del modelo multinomial se estiman utilizando máxima verosimilitud, similar a la regresión logística binaria.

47.53 Análisis de Supervivencia

El análisis de supervivencia se utiliza para modelar el tiempo hasta que ocurre un evento de interés, como la muerte o la falla de un componente.

47.53.1 Función de Supervivencia

La función de supervivencia $S(t)$ describe la probabilidad de que una observación sobreviva más allá del tiempo t :

$$S(t) = P(T > t)$$

47.53.2 Modelo de Riesgos Proporcionales de Cox

El modelo de Cox es un modelo de regresión semiparamétrico utilizado para analizar datos de supervivencia:

$$h(t|X) = h_0(t)e^{\beta_1X_1 + \dots + \beta_pX_p}$$

donde $h(t|X)$ es la tasa de riesgo en el tiempo t dado el vector de covariables X y $h_0(t)$ es la tasa de riesgo basal.

47.54 Implementación en R

47.54.1 Regresión Logística Multinomial

```
# Cargar el paquete necesario
library(nnet)

# Entrenar el modelo de regresión logística multinomial
model_multinom <- multinom(var1 ~ ., data = dataTrain)

# Resumen del modelo
summary(model_multinom)
```

47.54.2 Análisis de Supervivencia

```
# Cargar el paquete necesario
library(survival)

# Crear el objeto de supervivencia
surv_object <- Surv(time = data$time, event = data$status)

# Ajustar el modelo de Cox
model_cox <- coxph(surv_object ~ var1 + var2, data = data)

# Resumen del modelo
summary(model_cox)
```

Chapter 48

Numérico

48.1 Conceptos Básicos de la Regresión Logística

La regresión logística es una técnica de modelado estadístico utilizada para predecir la probabilidad de un evento binario (es decir, un evento que tiene dos posibles resultados) en función de una o más variables independientes. A diferencia de la regresión lineal, que se utiliza para predecir valores continuos, la regresión logística se usa cuando la variable dependiente es categórica.

48.2 Diferencias entre Regresión Lineal y Logística

48.2.1 Regresión Lineal

La regresión lineal busca modelar la relación entre una variable dependiente continua Y y una o más variables independientes X_1, X_2, \dots, X_n mediante una ecuación de la forma:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

donde $\beta_0, \beta_1, \dots, \beta_n$ son los coeficientes del modelo y ϵ es el término de error.

48.2.2 Regresión Logística

La regresión logística, en cambio, modela la probabilidad de que un evento ocurra (por ejemplo, éxito vs. fracaso) utilizando la función logística. La variable dependiente Y es binaria, tomando valores de 0 o 1. La ecuación de la regresión logística es:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

donde p es la probabilidad de que $Y = 1$. La función logística es:

$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n)}}$$

48.3 Casos de Uso de la Regresión Logística

La regresión logística se utiliza en una variedad de campos para problemas de clasificación binaria, tales como:

- **Medicina:** Predicción de la presencia o ausencia de una enfermedad.
- **Marketing:** Determinación de la probabilidad de que un cliente compre un producto.
- **Finanzas:** Evaluación del riesgo de crédito, es decir, si un cliente va a incumplir o no con un préstamo.
- **Seguridad:** Detección de fraudes o intrusiones.

48.4 Implementación Básica en R

Para implementar una regresión logística en R, primero es necesario instalar y cargar los paquetes necesarios. Aquí se muestra un ejemplo básico de implementación:

48.4.1 Instalación y Configuración de R y RStudio

- Descargue e instale R desde <https://cran.r-project.org/>.
- Descargue e instale RStudio desde <https://rstudio.com/products/rstudio/download/>.

48.4.2 Introducción Básica a R

- Sintaxis básica de R.
- Operaciones básicas: asignación, operaciones aritméticas, funciones básicas.

48.4.3 Ejemplo de Regresión Logística en R

```
# Instalación del paquete necesario
install.packages("stats")

# Carga del paquete
library(stats)

# Ejemplo de conjunto de datos
data <- data.frame(
  outcome = c(1, 0, 1, 0, 1, 1, 0, 1, 0, 0),
  predictor = c(2.3, 1.9, 3.1, 2.8, 3.6, 2.4, 2.1, 3.3, 2.2, 1.7)
)

# Ajuste del modelo de regresión logística
model <- glm(outcome ~ predictor, data = data, family = binomial)

# Resumen del modelo
summary(model)
```

En este ejemplo, se utiliza el conjunto de datos ‘data’ que contiene una variable de resultado binaria ‘outcome’ y una variable predictora continua ‘predictor’. El modelo de regresión logística se ajusta utilizando la función `glm` con la familia binomial.

48.5 Día 1: Regresión Logística

Implementación Básica en R

Para implementar una regresión logística en R, primero es necesario instalar y cargar los paquetes necesarios.

Instalación y Configuración de R y RStudio

- Descargue e instale R desde <https://cran.r-project.org/>. Siga las instrucciones para su sistema operativo (Windows, MacOS, Linux).
- Descargue e instale RStudio desde <https://rstudio.com/products/rstudio/download/>.

48.5.1 Ejemplo de Regresión Logística en R

A continuación, se muestra un ejemplo de cómo ajustar un modelo de regresión logística en R utilizando un conjunto de datos simulado. El ejemplo incluye la instalación del paquete necesario, la carga de datos, el ajuste del modelo, y la interpretación de los resultados.

```
# Instalación del paquete necesario
install.packages("stats")

# Carga del paquete
library(stats)

# Ejemplo de conjunto de datos
data <- data.frame(
  outcome = c(1, 0, 1, 0, 1, 1, 0, 1, 0, 0),
  predictor = c(2.3, 1.9, 3.1, 2.8, 3.6, 2.4, 2.1, 3.3, 2.2, 1.7)
)

# Ajuste del modelo de regresión logística
model <- glm(outcome ~ predictor, data = data, family = binomial)

# Resumen del modelo
summary(model)
```

En este ejemplo, se utiliza el conjunto de datos *data* que contiene una variable de resultado binaria *outcome* y una variable predictora continua *predictor*. El modelo de regresión logística se ajusta utilizando la función `glm` con la familia binomial. La función `summary(model)` proporciona un resumen del modelo ajustado, incluyendo los coeficientes estimados, sus errores estándar, valores z, y p-valores.

- **Coefficientes:** Los coeficientes estimados β_0 y β_1 indican la dirección y magnitud de la relación entre las variables predictoras y la probabilidad del resultado.

- **Errores Estándar:** Los errores estándar proporcionan una medida de la precisión de los coeficientes estimados.
- **Valores z y p-valores:** Los valores z y p-valores se utilizan para evaluar la significancia estadística de los coeficientes. Un p-valor pequeño (generalmente ≤ 0.05) indica que el coeficiente es significativamente diferente de cero.

Este es solo un ejemplo básico, en aplicaciones reales, es posible que necesites realizar más análisis y validaciones, como la evaluación de la bondad de ajuste del modelo, el diagnóstico de posibles problemas de multicolinealidad, y la validación cruzada del modelo.

```
# Archivo: regresionlogistica.R

# Instalación del paquete necesario
#install.packages("stats")

# Carga del paquete
library(stats)

# Fijar la semilla para reproducibilidad
set.seed(123)

# Número de observaciones
n <- 100

# Generar las variables independientes X1, X2, ..., X15
# Creamos una matriz de tamaño n x 15 con valores generados aleatoriamente de una
# distribución normal
X <- as.data.frame(matrix(rnorm(n * 15), nrow = n, ncol = 15))
colnames(X) <- paste0("X", 1:15) # Nombramos las columnas como X1, X2, ..., X15

# Coeficientes verdaderos para las variables independientes
# Generamos un vector de 16 coeficientes (incluyendo el intercepto) aleatorios entre -1 y 1
beta <- runif(16, -1, 1) # 15 coeficientes más el intercepto

# Generar el término lineal
# Calculamos el término lineal utilizando los coeficientes y las variables independientes
linear_term <- beta[1] + as.matrix(X) %*% beta[-1]

# Generar la probabilidad utilizando la función logística
# Calculamos las probabilidades utilizando la función logística
p <- 1 / (1 + exp(-linear_term))

# Generar la variable dependiente binaria Y
# Generamos valores binarios (0 o 1) utilizando las probabilidades calculadas
Y <- rbinom(n, 1, p)

# Combinar las variables independientes y la variable dependiente en un data frame
data <- cbind(Y, X)

# Dividir el conjunto de datos en entrenamiento y prueba
set.seed(123) # Fijar la semilla para reproducibilidad
```

```

train_indices <- sample(1:n, size = 0.7 * n) # 70% de los datos para entrenamiento
train_set <- data[train_indices, ] # Conjunto de entrenamiento
test_set <- data[-train_indices, ] # Conjunto de prueba

# Ajuste del modelo de regresión logística en el conjunto de entrenamiento
# Ajustamos un modelo de regresión logística utilizando las variables independientes
para predecir Y
model <- glm(Y ~ ., data = train_set, family = binomial)

# Resumen del modelo
# Mostramos un resumen del modelo ajustado
summary(model)

# Guardar el modelo y los resultados en un archivo
# Guardamos el modelo ajustado en un archivo .RData
save(model, file = "regresion_logistica_modelo.RData")

# Guardar los datos simulados en archivos CSV
# Guardamos los conjuntos de datos de entrenamiento y prueba en archivos CSV
write.csv(train_set, "datos_entrenamiento_regresion_logistica.csv", row.names = FALSE)
write.csv(test_set, "datos_prueba_regresion_logistica.csv", row.names = FALSE)

# Hacer predicciones en el conjunto de prueba
# Utilizamos el modelo ajustado para hacer predicciones en el conjunto de prueba
test_set$prob_pred <- predict(model, newdata = test_set, type = "response")
test_set$Y_pred <- ifelse(test_set$prob_pred > 0.5, 1, 0)
# Convertimos probabilidades a clases binarias

# Calcular la precisión de las predicciones
# Calculamos la precisión de las predicciones comparando con los valores reales de Y
accuracy <- mean(test_set$Y_pred == test_set$Y)
cat("La precisión del modelo en el conjunto de prueba es:", accuracy, "\n")

# Guardar las predicciones en un archivo CSV
# Guardamos las predicciones y las probabilidades predichas en un archivo CSV
write.csv(test_set, "predicciones_regresion_logistica.csv", row.names = FALSE)

# Graficar los coeficientes estimados
# Graficamos los coeficientes estimados del modelo ajustado
plot(coef(model), main = "Coeficientes Estimados del Modelo de Regresión Logística",
      xlab = "Variables", ylab = "Coeficientes", type = "h", col = "blue")
abline(h = 0, col = "red", lwd = 2)

# Mostrar un mensaje indicando que el proceso ha finalizado
cat("El modelo de regresión logística se ha ajustado, se han hecho predicciones y los resulta

```

48.5.2 Aplicación a Datos de Cáncer - Parte I

A continuación, se muestra un ejemplo de cómo ajustar un modelo de regresión logística en R utilizando el conjunto de datos del cáncer de mama de Wisconsin.

```

# Archivo: regresionlogistica_cancer.R

# Instalación del paquete necesario
install.packages("mlbench")
install.packages("dplyr")

# Carga de los paquetes
library(mlbench)
library(dplyr)

# Cargar el conjunto de datos BreastCancer
data("BreastCancer")

# Ver las primeras filas del conjunto de datos
head(BreastCancer)

# Preprocesamiento de los datos
# Eliminar la columna de identificación y filas con valores faltantes
breast_cancer_clean <- BreastCancer %>%
  select(-Id) %>%
  na.omit()

# Convertir la variable 'Class' a factor binario
breast_cancer_clean$Class <- ifelse(breast_cancer_clean$Class == "malignant", 1, 0)
breast_cancer_clean$Class <- as.factor(breast_cancer_clean$Class)

# Convertir las demás columnas a numéricas
breast_cancer_clean[, 1:9] <- lapply(breast_cancer_clean[, 1:9], as.numeric)

# Dividir el conjunto de datos en entrenamiento (70%) y prueba (30%)
set.seed(123)
train_indices <- sample(1:nrow(breast_cancer_clean), size = 0.7 * nrow(breast_cancer_clean))
train_set <- breast_cancer_clean[train_indices, ]
test_set <- breast_cancer_clean[-train_indices, ]

# Ajuste del modelo de regresión logística en el conjunto de entrenamiento
model <- glm(Class ~ ., data = train_set, family = binomial)

# Resumen del modelo
summary(model)

# Guardar el modelo y los resultados en un archivo
save(model, file = "regresion_logistica_cancer_modelo.RData")

# Guardar los datos simulados en archivos CSV
write.csv(train_set, "datos_entrenamiento_cancer.csv", row.names = FALSE)
write.csv(test_set, "datos_prueba_cancer.csv", row.names = FALSE)

# Hacer predicciones en el conjunto de prueba
test_set$prob_pred <- predict(model, newdata = test_set, type = "response")
test_set$Class_pred <- ifelse(test_set$prob_pred > 0.5, 1, 0)

```

```
# Calcular la precisión de las predicciones
accuracy <- mean(test_set$Class_pred == test_set$Class)
cat("La precisión del modelo en el conjunto de prueba es:", accuracy, "\n")

# Guardar las predicciones en un archivo CSV
write.csv(test_set, "predicciones_cancer.csv", row.names = FALSE)

# Graficar los coeficientes estimados
plot(coef(model), main = "Coeficientes Estimados del Modelo de Regresión Logística",
      xlab = "Variables", ylab = "Coeficientes", type = "h", col = "blue")
abline(h = 0, col = "red", lwd = 2)

# Mostrar un mensaje indicando que el proceso ha finalizado
cat("El modelo de regresión logística se ha ajustado, se han hecho predicciones y los resulta
```

Descripción del Código

Instalación y Carga de Paquetes:

Instalamos y cargamos el paquete `stats` necesario para la regresión logística.

Generación de Datos Simulados:

- Fijamos una semilla para la reproducibilidad.
- Generamos un conjunto de datos con 100 observaciones y 15 variables independientes (X_1 , X_2 , ..., X_{15}) usando una distribución normal.
- Definimos los coeficientes verdaderos para las variables independientes y calculamos el término lineal.
- Calculamos las probabilidades usando la función logística y generamos una variable dependiente binaria Y basada en esas probabilidades.
- Combinamos las variables independientes y la variable dependiente en un `data frame`.

División de Datos en Conjuntos de Entrenamiento y Prueba:

- Dividimos los datos en un conjunto de entrenamiento (70%) y un conjunto de prueba (30%).

Ajuste del Modelo de Regresión Logística:

- Ajustamos un modelo de regresión logística en el conjunto de entrenamiento.
- Mostramos un resumen del modelo ajustado.

Guardado de Datos y Modelo:

- Guardamos el modelo ajustado en un archivo `.RData`.
- Guardamos los conjuntos de datos de entrenamiento y prueba en archivos CSV.

Predicciones y Evaluación del Modelo:

- Hacemos predicciones en el conjunto de prueba utilizando el modelo ajustado.
- Calculamos la precisión de las predicciones comparando con los valores reales de Y .

- Guardamos las predicciones y las probabilidades predichas en un archivo CSV.

Visualización de los Coeficientes del Modelo:

- Graficamos los coeficientes estimados del modelo ajustado.
- Mostramos un mensaje indicando que el proceso ha finalizado.

Para ejecutar este script, guarda el código en un archivo llamado *regresionlogistica.R*, abre R o RStudio, navega hasta el directorio donde guardaste el archivo y ejecuta el script usando `source("regresionlogistica.R")`.

Ejemplo Titanic

Cuando realizas una regresión logística, obtienes coeficientes para cada variable independiente en tu modelo. Estos coeficientes indican la dirección y la magnitud de la relación entre cada variable independiente y la variable dependiente (en este caso, *Survived*).

Interpretación de los Coeficientes

- **Intercepto** (*(Intercept)*): Este coeficiente representa el logaritmo de las probabilidades (log-odds) de que *Survived* sea 1 (supervivencia) cuando todas las variables independientes son cero.
- **Pclass**: El coeficiente asociado con *Pclass* indica cómo cambia el log-odds de supervivencia con cada incremento en la clase del pasajero. Si el coeficiente es negativo, sugiere que una clase más alta (por ejemplo, de primera clase a tercera clase) reduce las probabilidades de supervivencia.
- **Sex**: Este coeficiente muestra el efecto de ser hombre o mujer en las probabilidades de supervivencia. Generalmente, se espera que el coeficiente sea positivo para *female* indicando que las mujeres tenían mayores probabilidades de sobrevivir.
- **Age**: El coeficiente de *Age* indica cómo cambia el log-odds de supervivencia con cada año de incremento en la edad. Un coeficiente negativo sugiere que la probabilidad de supervivencia disminuye con la edad.
- **SibSp** y **Parch**: Estos coeficientes indican el efecto del número de hermanos/cónyuges a bordo y padres/hijos a bordo en las probabilidades de supervivencia.
- **Fare**: Este coeficiente indica el efecto del precio del billete en las probabilidades de supervivencia. Un coeficiente positivo sugiere que pagar más por el billete se asocia con mayores probabilidades de supervivencia.

Estadísticas de Ajuste del Modelo

El resumen del modelo (`summary(model)`) incluye varias estadísticas importantes:

- **Estadísticos z y p-valores**: Estas estadísticas indican la significancia de cada coeficiente. Un p-valor bajo (generalmente ≤ 0.05) sugiere que la variable es un predictor significativo de la variable dependiente.
- **Desviación Residual**: La desviación residual mide la calidad del ajuste del modelo. Valores más bajos indican un mejor ajuste.
- **AIC (Akaike Information Criterion)**: El AIC es una medida de la calidad del modelo que toma en cuenta tanto la bondad del ajuste como la complejidad del modelo. Modelos con AIC más bajo son preferidos.

Precisión del Modelo

La precisión del modelo en el conjunto de prueba es una métrica importante para evaluar el rendimiento del modelo. La precisión se calcula como el número de predicciones correctas dividido por el número total de predicciones.

Ejemplo de Resultados

Supongamos que la precisión del modelo es 0.78 (78%). Esto significa que el modelo correctamente predijo el estado de supervivencia del 78% de los pasajeros en el conjunto de prueba.

Matriz de Confusión y Otras Métricas

Además de la precisión, otras métricas como la matriz de confusión, la sensibilidad, la especificidad, y el área bajo la curva ROC (AUC-ROC) también pueden proporcionar una visión más completa del rendimiento del modelo.

Matriz de Confusión

- **Verdaderos Positivos (TP)**: Número de pasajeros que sobrevivieron y fueron predichos como sobrevivientes.
- **Verdaderos Negativos (TN)**: Número de pasajeros que no sobrevivieron y fueron predichos como no sobrevivientes.
- **Falsos Positivos (FP)**: Número de pasajeros que no sobrevivieron pero fueron predichos como sobrevivientes.
- **Falsos Negativos (FN)**: Número de pasajeros que sobrevivieron pero fueron predichos como no sobrevivientes.

Ejemplo de Cálculo de Métricas

```
# Calcular la matriz de confusión
table(test_set$Survived, test_set$Survived_pred)

# Calcular sensibilidad y especificidad
sensitivity <- sum(test_set$Survived == 1 & test_set$Survived_pred == 1) / sum(test_set$Survived == 1)
specificity <- sum(test_set$Survived == 0 & test_set$Survived_pred == 0) / sum(test_set$Survived == 0)

# Calcular AUC-ROC
library(pROC)
roc_curve <- roc(test_set$Survived, test_set$prob_pred)
auc(roc_curve)
```

Visualización de Resultados

Graficar los coeficientes del modelo, la curva ROC y otras visualizaciones ayudan a entender mejor el rendimiento y la importancia de cada variable en el modelo.

```
# Graficar la curva ROC
plot(roc_curve, main = "Curva ROC para el Modelo de Regresión Logística")
```

Resumen Final

El modelo de regresión logística aplicado al conjunto de datos del Titanic proporciona una forma de entender cómo diferentes características de los pasajeros influyen en sus probabilidades de supervivencia. La interpretación de los coeficientes del modelo, las estadísticas de ajuste, y la precisión del modelo en el conjunto de prueba son fundamentales para evaluar el rendimiento y la utilidad del modelo en hacer predicciones sobre la supervivencia de los pasajeros del Titanic.

48.5.3 Simulación de Datos de Cáncer - Parte II

Aquí se presenta un ejemplo de cómo realizar una regresión logística utilizando datos simulados de pacientes con cáncer.

```
#---- Archivo: cancerLogRegSimulado.R ----

# Instalación del paquete necesario
if (!requireNamespace("dplyr", quietly = TRUE)) {
  install.packages("dplyr")
}

# Carga del paquete
library(dplyr)

# Fijar la semilla para reproducibilidad
set.seed(123)

# Número de observaciones
n <- 150

# Generar las variables independientes X1, X2, ..., X15
# Creamos una matriz de tamaño n x 15 con valores generados aleatoriamente de una
distribución normal
X <- as.data.frame(matrix(rnorm(n * 15), nrow = n, ncol = 15))
colnames(X) <- paste0("X", 1:15) # Nombramos las columnas como X1, X2, ..., X15

# Coeficientes verdaderos para las variables independientes
# Generamos un vector de 16 coeficientes (incluyendo el intercepto) aleatorios entre -1 y 1
beta <- runif(16, -1, 1) # 15 coeficientes más el intercepto

# Generar el término lineal
# Calculamos el término lineal utilizando los coeficientes y las variables independientes
linear_term <- beta[1] + as.matrix(X) %*% beta[-1]

# Generar la probabilidad utilizando la función logística
# Calculamos las probabilidades utilizando la función logística
p <- 1 / (1 + exp(-linear_term))

# Generar la variable dependiente binaria Y
# Generamos valores binarios (0 o 1) utilizando las probabilidades calculadas
Y <- rbinom(n, 1, p)
```

```

# Combinar las variables independientes y la variable dependiente en un data frame
data <- cbind(Y, X)

# Dividir el conjunto de datos en entrenamiento y prueba
set.seed(123) # Fijar la semilla para reproducibilidad
train_indices <- sample(1:n, size = 0.7 * n) # 70% de los datos para entrenamiento
train_set <- data[train_indices, ] # Conjunto de entrenamiento
test_set <- data[-train_indices, ] # Conjunto de prueba

# Ajuste del modelo de regresión logística en el conjunto de entrenamiento
# Ajustamos un modelo de regresión logística utilizando las variables independientes
para predecir Y
model <- glm(Y ~ ., data = train_set, family = binomial)

# Resumen del modelo
# Mostramos un resumen del modelo ajustado
summary(model)

# Guardar el modelo y los resultados en un archivo
# Guardamos el modelo ajustado en un archivo .RData
save(model, file = "regresion_logistica_cancer_modelo_simulado.RData")

# Guardar los datos simulados en archivos CSV
# Guardamos los conjuntos de datos de entrenamiento y prueba en archivos CSV
write.csv(train_set, "datos_entrenamiento_cancer_simulado.csv", row.names = FALSE)
write.csv(test_set, "datos_prueba_cancer_simulado.csv", row.names = FALSE)

# Hacer predicciones en el conjunto de prueba
# Utilizamos el modelo ajustado para hacer predicciones en el conjunto de prueba
test_set$prob_pred <- predict(model, newdata = test_set, type = "response")
test_set$Y_pred <- ifelse(test_set$prob_pred > 0.5, 1, 0)
# Convertimos probabilidades a clases binarias

# Calcular la precisión de las predicciones
# Calculamos la precisión de las predicciones comparando con los valores reales de Y
accuracy <- mean(test_set$Y_pred == test_set$Y)
cat("La precisión del modelo en el conjunto de prueba es:", accuracy, "\n")

# Guardar las predicciones en un archivo CSV
# Guardamos las predicciones y las probabilidades predichas en un archivo CSV
write.csv(test_set, "predicciones_cancer_simulado.csv", row.names = FALSE)

# Graficar los coeficientes estimados
# Graficamos los coeficientes estimados del modelo ajustado
plot(coef(model), main = "Coeficientes Estimados del Modelo de Regresión Logística",
      xlab = "Variables", ylab = "Coeficientes", type = "h", col = "blue")
abline(h = 0, col = "red", lwd = 2)

# Mostrar un mensaje indicando que el proceso ha finalizado
cat("El modelo de regresión logística se ha ajustado, se han hecho predicciones
y los resultados se han guardado en 'regresion_logistica_cancer_modelo_simulado.RData'.\n")

```

48.5.4 Simulación de Datos de Cáncer - Parte III

En un estudio sobre cáncer, especialmente en el contexto del cáncer de mama, las principales mediciones suelen incluir una variedad de características clínicas y patológicas. Aquí hay algunas de las principales mediciones que se tienen en cuenta:

- **Tamaño del Tumor:** Medición del diámetro del tumor.
- **Estado de los Ganglios Linfáticos:** Número de ganglios linfáticos afectados.
- **Grado del Tumor:** Clasificación del tumor basada en la apariencia de las células cancerosas.
- **Receptores Hormonales:** Estado de los receptores de estrógeno y progesterona.
- **Estado HER2:** Expresión del receptor 2 del factor de crecimiento epidérmico humano.
- **Ki-67:** Índice de proliferación celular.
- **Edad del Paciente:** Edad en el momento del diagnóstico.
- **Histopatología:** Tipo y subtipo histológico del cáncer.
- **Márgenes Quirúrgicos:** Estado de los márgenes después de la cirugía (si están libres de cáncer o no).
- **Invasión Linfovascular:** Presencia de células cancerosas en los vasos linfáticos o sanguíneos.
- **Tratamientos Previos:** Tipos de tratamientos recibidos antes del diagnóstico (quimioterapia, radioterapia, etc.).
- **Tipo de Cirugía:** Tipo de procedimiento quirúrgico realizado (mastectomía, lumpectomía, etc.).
- **Metástasis:** Presencia de metástasis y ubicación de las mismas.
- **Índice de Masa Corporal (IMC):** Relación entre el peso y la altura del paciente.
- **Marcadores Genéticos:** Presencia de mutaciones genéticas específicas (BRCA1, BRCA2, etc.).

Estas mediciones proporcionan una visión integral del estado del cáncer y se utilizan para planificar el tratamiento y predecir el pronóstico.

A continuación, se muestra un ejemplo de cómo ajustar un modelo de regresión logística en R utilizando un conjunto de datos simulado con estas mediciones.

```
# Archivo: simulcorrectedCancer.R

# Instalación del paquete necesario
if (!requireNamespace("dplyr", quietly = TRUE)) {
  install.packages("dplyr")
}

# Carga del paquete
library(dplyr)

# Fijar la semilla para reproducibilidad
set.seed(123)

# Número de observaciones
n <- 1500

# Simulación de las variables independientes
```

```

# Tamaño del Tumor (en cm)
Tumor_Size <- rnorm(n, mean = 3, sd = 1.5)

# Estado de los Ganglios Linfáticos (número de ganglios afectados)
Lymph_Nodes <- rpois(n, lambda = 3)

# Grado del Tumor (1 a 3)
Tumor_Grade <- sample(1:3, n, replace = TRUE)

# Receptores Hormonales (0: negativo, 1: positivo)
Estrogen_Receptor <- rbinom(n, 1, 0.7)
Progesterone_Receptor <- rbinom(n, 1, 0.7)

# Estado HER2 (0: negativo, 1: positivo)
HER2_Status <- rbinom(n, 1, 0.3)

# Ki-67 (% de células proliferativas)
Ki_67 <- rnorm(n, mean = 20, sd = 10)

# Edad del Paciente (años)
Age <- rnorm(n, mean = 50, sd = 10)

# Histopatología (1: ductal, 2: lobular, 3: otros)
Histopathology <- sample(1:3, n, replace = TRUE)

# Márgenes Quirúrgicos (0: positivo, 1: negativo)
Surgical_Margins <- rbinom(n, 1, 0.8)

# Invasión Linfovascular (0: no, 1: sí)
Lymphovascular_Invasion <- rbinom(n, 1, 0.4)

# Tratamientos Previos (0: no, 1: sí)
Prior_Treatments <- rbinom(n, 1, 0.5)

# Tipo de Cirugía (0: mastectomía, 1: lumpectomía)
Surgery_Type <- rbinom(n, 1, 0.5)

# Metástasis (0: no, 1: sí)
Metastasis <- rbinom(n, 1, 0.2)

# Índice de Masa Corporal (IMC)
BMI <- rnorm(n, mean = 25, sd = 5)

# Marcadores Genéticos (0: negativo, 1: positivo)
Genetic_Markers <- rbinom(n, 1, 0.1)

# Generar la variable dependiente binaria Y (sobrevivencia 0: no, 1: sí)
# Utilizaremos una combinación arbitraria de las variables para generar Y
linear_term <- -1 + 0.5 * Tumor_Size - 0.3 * Lymph_Nodes + 0.2 * Tumor_Grade +
  0.4 * Estrogen_Receptor + 0.3 * Progesterone_Receptor - 0.2 * HER2_Status +
  0.1 * Ki_67 - 0.05 * Age + 0.3 * Surgical_Margins - 0.4 * Lymphovascular_Invasion +

```

```

0.2 * Prior_Treatments + 0.1 * Surgery_Type - 0.5 * Metastasis + 0.01 * BMI +
0.2 * Genetic_Markers
p <- 1 / (1 + exp(-linear_term))
Y <- rbinom(n, 1, p)

# Combinar las variables independientes y la variable dependiente en un data frame
data <- data.frame(Y, Tumor_Size, Lymph_Nodes, Tumor_Grade, Estrogen_Receptor,
                  Progesterone_Receptor, HER2_Status, Ki_67, Age, Histopathology,
                  Surgical_Margins, Lymphovascular_Invasion, Prior_Treatments,
                  Surgery_Type, Metastasis, BMI, Genetic_Markers)

# Dividir el conjunto de datos en entrenamiento y prueba
set.seed(123) # Fijar la semilla para reproducibilidad
train_indices <- sample(1:n, size = 0.7 * n) # 70% de los datos para entrenamiento
train_set <- data[train_indices, ] # Conjunto de entrenamiento
test_set <- data[-train_indices, ] # Conjunto de prueba

# Ajuste del modelo de regresión logística en el conjunto de entrenamiento
# Ajustamos un modelo de regresión logística utilizando las variables independientes para
predecir Y
model <- glm(Y ~ ., data = train_set, family = binomial)

# Resumen del modelo
# Mostramos un resumen del modelo ajustado
summary(model)

# Guardar el modelo y los resultados en un archivo
# Guardamos el modelo ajustado en un archivo .RData
save(model, file = "regresion_logistica_cancer_modelo_simulado.RData")

# Guardar los datos simulados en archivos CSV
# Guardamos los conjuntos de datos de entrenamiento y prueba en archivos CSV
write.csv(train_set, "datos_entrenamiento_cancer_simulado.csv", row.names = FALSE)
write.csv(test_set, "datos_prueba_cancer_simulado.csv", row.names = FALSE)

# Hacer predicciones en el conjunto de prueba
# Utilizamos el modelo ajustado para hacer predicciones en el conjunto de prueba
test_set$prob_pred <- predict(model, newdata = test_set, type = "response")
test_set$Y_pred <- ifelse(test_set$prob_pred > 0.5, 1, 0)
# Convertimos probabilidades a clases binarias

# Calcular la precisión de las predicciones
# Calculamos la precisión de las predicciones comparando con los valores reales de Y
accuracy <- mean(test_set$Y_pred == test_set$Y)
cat("La precisión del modelo en el conjunto de prueba es:", accuracy, "\n")

# Guardar las predicciones en un archivo CSV
# Guardamos las predicciones y las probabilidades predichas en un archivo CSV
write.csv(test_set, "predicciones_cancer_simulado.csv", row.names = FALSE)

# Graficar los coeficientes estimados

```

```
# Graficamos los coeficientes estimados del modelo ajustado
plot(coef(model), main = "Coeficientes Estimados del Modelo de Regresión Logística",
      xlab = "Variables", ylab = "Coeficientes", type = "h", col = "blue")
abline(h = 0, col = "red", lwd = 2)

# Mostrar un mensaje indicando que el proceso ha finalizado
cat("El modelo de regresión logística se ha ajustado, se han hecho predicciones
y los resultados se han guardado en 'regresion_logistica_cancer_modelo_simulado.RData'.\n")
```


Bibliography

- [1] Darlington RB. *Regression and Linear Models*. Columbus, OH: McGraw-Hill Publishing Company, 1990.
- [2] Tabachnick BG, Fidell LS. *Using Multivariate Statistics*. 5th ed. Boston, MA: Pearson Education, Inc., 2007.
- [3] Hosmer DW, Lemeshow SL. *Applied Logistic Regression*. 2nd ed. Hoboken, NJ: Wiley-Interscience, 2000.
- [4] Campbell DT, Stanley JC. *Experimental and Quasi-experimental Designs for Research*. Boston, MA: Houghton Mifflin Co., 1963.
- [5] Stokes ME, Davis CS, Koch GG. *Categorical Data Analysis Using the SAS System*. 2nd ed. Cary, NC: SAS Institute, Inc., 2000.
- [6] Newgard CD, Hedges JR, Arthur M, Mullins RJ. Advanced statistics: the propensity score—a method for estimating treatment effect in observational research. *Acad Emerg Med*. 2004; **11**:953–961.
- [7] Newgard CD, Haukoos JS. Advanced statistics: missing data in clinical research—part 2: multiple imputation. *Acad Emerg Med*. 2007; **14**:669–678.
- [8] Allison PD. *Logistic Regression Using the SAS System: Theory and Application*. Cary, NC: SAS Institute, Inc., 1999.
- [9] Peduzzi P, Concato J, Kemper E, Holford TR, Feinstein AR. A simulation study of the number of events per variable in logistic regression analysis. *J Clin Epidemiol*. 1996; **49**:1373–1379.
- [10] Agresti A. *An Introduction to Categorical Data Analysis*. Hoboken, NJ: Wiley, 2007.
- [11] Feinstein AR. *Multivariable Analysis: An Introduction*. New Haven, CT: Yale University Press, 1996.
- [12] Altman DG, Royston P. What Do We Mean by Validating a Prognostic Model? *Stats Med*. 2000; **19**:453–473.
- [13] Kohavi R. A study of cross-validation and bootstrap for accuracy estimation and model selection. In: *Proceedings of the 14th International Joint Conference on Artificial Intelligence (IJCAI)*. Montreal, Quebec, Canada, August 20–25, 1995. 1995:1137–1143.
- [14] Efron B, Tibshirani R. *An Introduction to the Bootstrap*. New York: Chapman & Hall, 1993.
- [15] Miller ME, Hiu SL, Tierney WM. Validation techniques for logistic regression models. *Stat Med*. 1991; **10**:1213–1226.
- [16] Hosmer DW, Hosmer T, Le Cessie S, Lemeshow S. A comparison of goodness-of-fit tests for the logistic regression model. *Stat Med*. 1997; **16**:965–980.
- [17] Kuss O. Global goodness-of-fit tests in logistic regression with sparse data. *Stat Med*. 2002; **21**:3789–3801.
- [18] Zou KH, O'Malley AJ, Mauri L. Receiver-operating characteristic analysis for evaluating diagnostic tests and predictive models. *Circulation*. 2007; **115**:654–657.

- [19] Mazurenko, S., Prokop, Z., and Damborsky, J. (2019). Machine learning in enzyme engineering. *ACS Catalysis*, 10(2), 1210-1223.
- [20] Yang, Y.; Niroula, A.; Shen, B.; Vihinen, M. PON-Sol: Prediction of Effects of Amino Acid Substitutions on Protein Solubility. *Bioinformatics* 2016, 32, 2032-2034.
- [21] Folkman, L.; Stantic, B.; Sattar, A.; Zhou, Y. EASE-MM: Sequence-Based Prediction of Mutation-Induced Stability Changes with Feature-Based Multiple Models. *J. Mol. Biol.* 2016, 428, 1394-1405.
- [22] Teng, S.; Srivastava, A. K.; Wang, L. Sequence Feature-Based Prediction of Protein Stability Changes upon Amino Acid Substitutions. *BMC Genomics* 2010, 11, S5.
- [23] Huang, L.; Gromiha, M. M.; Ho, S. iPTREE-STAB: Interpretable Decision Tree Based Method for Predicting Protein Stability Changes Upon Mutations. *Bioinformatics* 2007, 23, 1292-1293.
- [24] Koskinen, P.; Toronen, P.; Nokso-Koivisto, J.; Holm, L. PANNZER: High-Throughput Functional Annotation of Uncharacterized Proteins in an Error-Prone Environment. *Bioinformatics* 2015, 31, 1544-1552.
- [25] De Ferrari, L.; Mitchell, J. B. From Sequence to Enzyme Mechanism Using Multi-Label Machine Learning. *BMC Bioinf.* 2014, 15, 150.
- [26] Falda, M.; Toppo, S.; Pescarolo, A.; Lavezzo, E.; Di Camillo, B.; Facchinetti, A.; Cilia, E.; Velasco, R.; Fontana, P. Argot2: A Large Scale Function Prediction Tool Relying on Semantic Similarity of Weighted Gene Ontology Terms. *BMC Bioinf.* 2012, 13, S14.
- [27] Cozzetto, D.; Buchan, D. W.; Bryson, K.; Jones, D. T. Protein Function Prediction by Massive Integration of Evolutionary Analyses and Multiple Data Sources. *BMC Bioinf.* 2013, 14, S1.
- [28] Kulski, J. Next Generation Sequencing: Advances, Applications and Challenges; InTechOpen: London, 2016.
- [29] Straiton, J.; Free, T.; Sawyer, A.; Martin, J. From Sanger Sequencing to Genome Databases and Beyond. *BioTechniques* 2019, 66, 60-63.
- [30] Ardui, S.; Ameer, A.; Vermeesch, J. R.; Hestand, M. S. Single Molecule Real-Time (SMRT) Sequencing Comes of Age: Applications and Utilities for Medical Diagnostics. *Nucleic Acids Res.* 2018, 46, 2159-2168.
- [31] Kono, N., and Arakawa, K. (2019). Nanopore sequencing: Review of potential applications in functional genomics. *Development, growth and differentiation*, 61(5), 316-326.
- [32] Bunzel, H. A., Garrabou, X., Pott, M., and Hilvert, D. (2018). Speeding up enzyme discovery and engineering with ultrahigh-throughput methods. *Current opinion in structural biology*, 48, 149-156.
- [33] Wrenbeck, E. E., Faber, M. S., and Whitehead, T. A. (2017). Deep sequencing methods for protein engineering and design. *Current opinion in structural biology*, 45, 36-44.
- [34] Fowler, D. M., and Fields, S. (2014). Deep mutational scanning: a new style of protein science. *Nature methods*, 11(8), 801-807.
- [35] Gupta, K., and Varadarajan, R. (2018). Insights into protein structure, stability and function from saturation mutagenesis. *Current opinion in structural biology*, 50, 117-125.
- [36] UniProt Consortium. UniProt: A Worldwide Hub of Protein Knowledge. *Nucleic Acids Res.* 2018, 47, D506-D515.
- [37] Evans, R.; Jumper, J.; Kirkpatrick, J.; Sifre, L.; Green, T.; Qin, C.; Zidek, A.; Nelson, A.; Bridgland, A.; Penedones, H.; Petersen, S.; Simonyan, K.; Jones, D. T.; Silver, D.; Kavukcuoglu, K.; Hassabis, D.; Senior, A. W. De Novo Structure Prediction with Deep Learning Based Scoring. In *Thirteenth Critical Assessment of Techniques for Protein Structure Prediction Abstracts*; 2018; pp 11-12.
- [38] Kinch, L. N.; Shi, S.; Cheng, H.; Cong, Q.; Pei, J.; Mariani, V.; Schwede, T.; Grishin, N. V. CASP9 Target Classification. *Proteins: Struct., Funct., Genet.* 2011, 79, 21-36.

- [39] Shehu, A.; Barbará, D.; Molloy, K. A Survey of Computational Methods for Protein Function Prediction. In *Big Data Analytics in Genomics*; Wong, K. C., Ed.; Springer: Cham, 2016; pp 225-298.
- [40] Zhang, C.; Freddolino, P. L.; Zhang, Y. COFACTOR: Improved Protein Function Prediction by Combining Structure, Sequence and Protein-Protein Interaction Information. *Nucleic Acids Res.* 2017, 45, W291-W299.
- [41] Kumar, N.; Skolnick, J. EFICAz2. 5: Application of a High-Precision Enzyme Function Predictor to 396 Proteomes. *Bioinformatics* 2012, 28, 2687-2688.
- [42] Li, Y.; Wang, S.; Umarov, R.; Xie, B.; Fan, M.; Li, L.; Gao, X. DEEPRe: Sequence-Based Enzyme EC Number Prediction by Deep Learning. *Bioinformatics* 2018, 34, 760-769.
- [43] Yang, M.; Fehl, C.; Lees, K. V.; Lim, E. K.; Offen, W. A.; Davies, G. J.; Bowles, D. J.; Davidson, M. G.; Roberts, S. J.; Davis, B. G. Functional and Informatics Analysis Enables Glycosyltransferase Activity Prediction. *Nat. Chem. Biol.* 2018, 14, 1109-1117.
- [44] Niwa, T.; Ying, B. W.; Saito, K.; Jin, W.; Takada, S.; Ueda, T.; Taguchi, H. Bimodal Protein Solubility Distribution Revealed by an Aggregation Analysis of the Entire Ensemble of Escherichia Coli Proteins. *Proc. Natl. Acad. Sci. U. S. A.* 2009, 106, 4201-4206.
- [45] Klesmith, J. R.; Bacik, J. P.; Wrenbeck, E. E.; Michalczyk, R.; Whitehead, T. A. Trade-Offs Between Enzyme Fitness and Solubility Illuminated by Deep Mutational Scanning. *Proc. Natl. Acad. Sci. U. S. A.* 2017, 114, 2265-2270.
- [46] Ruiz-Blanco, Y. B.; Paz, W.; Green, J.; Marrero-Ponce, Y. ProtDCal: A Program to Compute General-Purpose-Numerical Descriptors for Sequences and 3D-Structures of Proteins. *BMC Bioinf.* 2015, 16, 162.
- [47] Han, X.; Wang, X.; Zhou, K. Develop Machine Learning-Based Regression Predictive Models for Engineering Protein Solubility. *Bioinformatics* 2019, 35, 4640-4646.
- [48] Musil, M.; Konegger, H.; Hon, J.; Bednar, D.; Damborsky, J. Computational Design of Stable and Soluble Biocatalysts. *ACS Catal.* 2019, 9, 1033-1054.
- [49] Li, G.; Dong, Y.; Reetz, M. T. Can Machine Learning Revolutionize Directed Evolution of Selective Enzymes? *Adv. Synth. Catal.* 2019, 361, 2377-2386.
- [50] Wu, Z.; Kan, S. B. J.; Lewis, R. D.; Wittmann, B. J.; Arnold, F. H. Machine Learning-Assisted Directed Protein Evolution with Combinatorial Libraries. *Proc. Natl. Acad. Sci. U. S. A.* 2019, 116, 8852-8858.
- [51] Wolpert, D. H.; Macready, W. G. No Free Lunch Theorems for Optimization. *IEEE Trans. Evol. Comput.* 1997, 1, 67-82.
- [52] Wolpert, D. H. The Lack of a Priori Distinctions between Learning Algorithms. *Neural Comput.* 1996, 8, 1341-1390.
- [53] Walsh, I.; Pollastri, G.; Tosatto, S. C. Correct Machine Learning on Protein Sequences: A Peer-Reviewing Perspective. *Briefings Bioinf.* 2016, 17, 831-840.
- [54] Rao, R.; Bhattacharya, N.; Thomas, N.; Duan, Y.; Chen, X.; Canny, J.; Abbeel, P.; Song, Y. S. Evaluating Protein Transfer Learning with TAPE. *arXiv preprint arXiv:1906.08230*, 2019.
- [55] Romero, P. A.; Krause, A.; Arnold, F. H. Navigating the Protein Fitness Landscape with Gaussian Processes. *Proc. Natl. Acad. Sci. U. S. A.* 2013, 110, E193-E201.
- [56] Eraslan, G.; Avsec, Z.; Gagneur, J.; Theis, F. J. Deep Learning: New Computational Modelling Techniques for Genomics. *Nat. Rev. Genet.* 2019, 20, 389-403.
- [57] Repecka, D.; Jauniskis, V.; Karpus, L.; Rembeza, E.; Zrimec, J.; Poviloniene, S.; Rokaitis, I.; Laurynenas, A.; Abuajwa, W.; Savolainen, O.; Meskys, R.; Engqvist, M. K. M.; Zelezniak, A. Expanding Functional Protein Sequence Space Using Generative Adversarial Networks. *bioRxiv* 2019, DOI: 10.1101/789719.

- [58] Riesselman, A. J.; Ingraham, J. B.; Marks, D. S. Deep Generative Models of Genetic Variation Capture the Effects of Mutations. *Nat. Methods* 2018, 15, 816-822.
- [59] Thornton, C.; Hutter, F.; Hoos, H. H.; Leyton-Brown, K. Auto-WEKA: Combined Selection and Hyperparameter Optimization of Classification Algorithms. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*; 2013; pp 847-855.
- [60] Polikar, R. Ensemble based systems in decision making. *IEEE Circuits and systems magazine* 2006, 6, 21-45.
- [61] Gammernan, A.; Vovk, V. Hedging Predictions in Machine Learning. *Comput. J.* 2007, 50, 151-163.
- [62] Samek, W.; Wiegand, T.; Müller, K. Explainable Artificial Intelligence: Understanding, Visualizing and Interpreting Deep Learning Models. *ITU Journal: ICT Discoveries* 2017, 39-48.
- [63] Shrikumar, A.; Greenside, P.; Kundaje, A. Learning Important Features through Propagating Activation differences. In *Proceedings of the 34th International Conference on Machine Learning*; 2017; Vol. 70, pp 3145-3153.
- [64] Simonyan, K.; Vedaldi, A.; Zisserman, A. Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps. *arXiv preprint arXiv:1312.6034* 2013.
- [65] Brookes, D. H.; Park, H.; Listgarten, J. Conditioning by Adaptive Sampling for Robust Design. In *Proceedings of the 36th International Conference on Machine Learning*; 2019; Vol. 97, pp 773-782.
- [66] Ribeiro, M. T.; Singh, S.; Guestrin, C. “Why Should I Trust You?” Explaining the Predictions of Any Classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*; 2016; pp 1135-1144.
- [67] Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; Fergus, R. Intriguing Properties of Neural Networks. *arXiv preprint arXiv:1312.6199* 201
- [68] Yu, M. K.; Ma, J.; Fisher, J.; Kreisberg, J. F.; Raphael, B. J.; Ideker, T. Visible Machine Learning for Biomedicine. *Cell* 2018, 173, 1562-1565.
- [69] Copley, S. D. Shining a light on enzyme promiscuity. *Curr. Opin. Struct. Biol.* 47, 167–175 (2017).
- [70] Nobeli, I., Favia, A. D. and Thornton, J. M. Protein promiscuity and its implications for biotechnology. *Nat. Biotechnol.* 27, 157–167 (2009)
- [71] Adrio, J. L. and Demain, A. L. Microbial enzymes: tools for biotechnological processes. *Biomolecules* 4, 117–139 (2014).
- [72] Wang, S. et al. Engineering a synthetic pathway for gentisate in *pseudomonas chlororaphis* p3. *Front. Bioeng. Biotechnol.* 8, 1588 (2021).
- [73] Wu, M.-C., Law, B., Wilkinson, B. and micklefield, J. Bioengineering natural product biosynthetic pathways for therapeutic applications. *Curr. Opin. Biotechnol.* 23, 931–940 (2012)
- [74] Rembeza, E., Boverio, A., Fraaije, M. W. and Engqvist, M. K. Discovery of two novel oxidases using a high-throughput activity screen. *ChemBioChem* 23, e202100510 (2022).
- [75] Longwell, C. K., Labanieh, L. and Cochran, J. R. High-throughput screening technologies for enzyme engineering. *Curr. Opin. Biotechnol.* 48, 196–202 (2017).
- [76] Black, G. W. et al. A high-throughput screening method for determining the substrate scope of nitrilases. *Chem. Commun.* 51, 2660–2662 (2015).
- [77] Pertusi, D. A. et al. Predicting novel substrates for enzymes with minimal experimental effort with active learning. *Metab. Eng.* 44, 171-181 (2017).
- [78] Mou, Z. et al. Machine learning-based prediction of enzyme substrate scope: Application to bacterial nitrilases. *Proteins Struct. Funct. Bioinf.* 89, 336-347 (2021).
- [79] Yang, M. et al. Functional and informatics analysis enables glycosyltransferase activity prediction. *Nat. Chem. Biol.* 14, 1109–1117 (2018).

- [80] Rottig, M., Rausch, C. and Kohlbacher, O. Combining structure and sequence information allows automated prediction of substrate specificities within enzyme families. *PLoS Comput. Biol.* 6, e1000636 (2010).
- [81] Chevrette, M. G., Aicheler, F., Kohlbacher, O., Currie, C. R. and Medema, M. H. Sandpuma: ensemble predictions of nonribosomal peptide chemistry reveal biosynthetic diversity across actinobacteria. *Bioinformatics* 33, 3202-3210 (2017).
- [82] Goldman, S., Das, R., Yang, K. K. and Coley, C. W. Machine learning modeling of family wide enzyme-substrate specificity screens. *PLoS Comput. Biol.* 18, e1009853 (2022).
- [83] Visani, G. M., Hughes, M. C. and Hassoun, S. Enzyme promiscuity prediction using hierarchy-informed multi-label classification *Bioinformatics* 37, 2017-2024 (2021).
- [84] Ryu, J. Y., Kim, H. U. and Lee, S. Y. Deep learning enables high-quality and high-throughput prediction of enzyme commission numbers. *PNAS* 116, 13996-14001 (2019).
- [85] Li, Y. et al. DEEPre: sequence-based enzyme EC number prediction by deep learning. *Bioinformatics* 34, 760-769 (2017).
- [86] Sanderson, T., Bileschi, M. L., Belanger, D. and Colwell, L. J. Proteinfer, deep neural networks for protein functional inference. *eLife* 12, e80942 (2023).
- [87] Bileschi, M. L. et al. Using deep learning to annotate the protein universe. *Nat. Biotechnol.* <https://doi.org/10.1038/s41587-021-01179-w> (2022).
- [88] Rembeza, E. and Engqvist, M. K. Experimental and computational investigation of enzyme functional annotations uncovers misannotation in the ec 1.1. 3.15 enzyme class. *PLoS Comput. Biol.* 17, e1009446 (2021).
- [89] Ozturk, H., Ozgur, A. and Ozkirimli, E. Deepdta: deep drugtarget binding affinity prediction. *Bioinformatics* 34, i821-i829 (2018).
- [90] Feng, Q., Dueva, E., Cherkasov, A. and Ester, M. Padme: A deep learning-based framework for drug-target interaction prediction. Preprint at <https://doi.org/10.48550/arXiv.1807.09741> (2018).
- [91] Karimi, M., Wu, D., Wang, Z. and Shen, Y. Deep affinity: interpretable deep learning of compound-protein affinity through UNIFIED recurrent and convolutional neural networks. *Bioinformatics* 35, 3329-3338 (2019).
- [92] Kroll, A., Engqvist, M. K., Heckmann, D. and Lercher, M. J. Deep learning allows genome-scale prediction of michaelis constants from structural features. *PLoS Biol.* 19, e3001402 (2021).
- [93] Li, F. et al. Deep learning-based k cat prediction enables improved enzyme-constrained model reconstruction. *Nat. Catal.* 5, 662-672 (2022).
- [94] Weininger, D. SMILES, a chemical language and information system. 1. introduction to methodology and encoding rules. *J. Chem. Inf. Comput. Sci.* 28, 31-36 (1988).
- [95] Rogers, D. and Hahn, M. Extended-connectivity fingerprints. *J. Chem. Inf. Model.* 50, 742-754 (2010).
- [96] Zhou, J. et al. Graph neural networks: A review of methods and applications. *AI Open* 1, 57-81 (2020).
- [97] Yang, K. et al. Analyzing learned molecular representations for property prediction. *J. Chem. Inf. Model.* 59, 3370-3388 (2019).
- [98] Rives, A. et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *PNAS* 118, e2016239118 (2021).
- [99] Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. and Church, G. M. Unified rational protein engineering with sequence-based deep representation learning. *Nat. Methods.* 16, 1315-1322 (2019).

- [100] Xu, Y. et al. Deep dive into machine learning models for protein engineering. *J. Chem. Inf. Model.* 60, 2773–2790 (2020).
- [101] Bekker, J. and Davis, J. Learning from positive and unlabeled data: A survey. *Mach. Learn.* 109, 719–760 (2020)
- [102] Kearnes, S., McCloskey, K., Berndl, M., Pande, V. and Riley, P. Molecular graph convolutions: moving beyond fingerprints. *J. Comput. -Aided Mol. Des.* 30, 595–608 (2016).
- [103] Duvenaud, D. K. et al. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems*, 2224–2232 (2015).
- [104] Zhou, J. et al. Graph neural networks: A review of methods and applications. *AI Open* 1, 57–81 (2020).
- [105] Hu, W. et al. Strategies for pre-training graph neural networks. Preprint at <https://doi.org/10.48550/arXiv.1905.12265> (2019).
- [106] Capela, F., Nouchi, V., Van Deursen, R., Tetko, I. V. and Godin, G. Multitask learning on graph neural networks applied to molecular property predictions. Preprint at <https://doi.org/10.48550/arXiv.1910.13124> (2019).
- [107] Vaswani, A. et al. Attention is all you need. In *Advances in neural information processing systems*, 5998–6008 (2017).
- [108] Suzeck, B. E. et al. Uniref clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics* 31, 926–932 (2015).
- [109] Elnaggar, A. et al. Prottrans: Towards cracking the language of life's code through self-supervised deep learning and high performance computing. *IEEE Trans. Pattern Anal. Mach. Intell.* PP <https://doi.org/10.1109/TPAMI.2021.3095381> (2021).
- [110] Wittmann, B. J., Johnston, K. E., Wu, Z., and Arnold, F. H. (2021). Advances in machine learning for directed evolution. *Current opinion in structural biology*, 69, 11–18.