

análisis de datos categóricos – categorical data analysis

Authored by
memjavad

November 12, 2025

RECOMMENDED CITATION

memjavad (2025). *análisis de datos categóricos – categorical data analysis*. Spanish Psychological Databases. Retrieved from <https://spanish.arabpsychology.com/?p=4061>

Análisis de Datos Categóricos

Primary Disciplinary Field(s): Estadística, Modelado Estadístico, Biometría, Ciencias Sociales Cuantitativas

1. Definición Fundamental y Naturaleza de los Datos Categóricos

El **Análisis de Datos Categóricos** (ADC) constituye una rama esencial de la estadística dedicada al estudio de variables que solo pueden tomar un número limitado de valores o categorías, en oposición a las variables continuas que pueden asumir cualquier valor dentro de un rango determinado. Este campo se ocupa de desarrollar y aplicar métodos robustos para describir la relación y asociación entre estas variables discretas. La necesidad de un conjunto de herramientas estadísticas especializado surge del hecho fundamental de que las variables categóricas no cumplen con los supuestos de normalidad y homocedasticidad que subyacen a gran parte de la estadística paramétrica tradicional, como la regresión lineal o las pruebas T. Por lo tanto, el ADC emplea técnicas que se basan en la distribución de frecuencias y las probabilidades, haciendo uso intensivo de modelos logarítmicos y distribuciones discretas, como la binomial o la de Poisson.

La distinción crucial en el ADC reside en la escala de medición de las variables. Cuando las variables son categóricas, la información intrínseca se relaciona con la pertenencia a un grupo y no con una magnitud medible en un sentido métrico. Esto implica que las operaciones aritméticas estándar (suma, resta, promediado) no tienen un significado estadístico válido o interpretativo directo. En su lugar, el análisis se centra en la proporción de observaciones que caen en cada categoría y en cómo esta distribución de proporciones varía en función de otras variables categóricas o continuas. La meta principal es determinar si existe una dependencia o asociación estadísticamente significativa entre las variables, y cuantificar la fuerza de esa relación mediante medidas como la razón de probabilidades (*odds ratio*).

Este enfoque analítico no solo es descriptivo, sino también predictivo. Aunque históricamente el campo se centró en las pruebas de independencia, las metodologías modernas del ADC han evolucionado para incluir potentes herramientas de modelado predictivo. Estas herramientas permiten a los investigadores predecir la probabilidad de que un resultado caiga en una categoría específica, basándose en un conjunto de variables predictoras. Por ejemplo, en epidemiología, el ADC es fundamental para predecir la probabilidad de desarrollar una enfermedad (resultado binario) en función de factores de riesgo categóricos (como el tabaquismo o la exposición ambiental). La comprensión profunda de la naturaleza discreta de los datos es, por lo tanto, el punto de partida para cualquier análisis categórico riguroso.

2. Tipologías de Datos Categóricos y Escalas de Medición

Dentro del espectro de variables categóricas, es imperativo diferenciar entre sus principales subtipos, ya que la elección del método estadístico adecuado depende directamente de la escala de medición. El primer subtipo es la **variable nominal**, donde las categorías son etiquetas que no poseen ningún orden o jerarquía intrínseca. Ejemplos clásicos incluyen la nacionalidad, el género o el color de ojos. Aunque se les puede asignar un código numérico para facilitar la entrada de datos, estos números son meros identificadores; no se puede afirmar que la categoría '1' sea mayor o menor que la categoría '2'. El análisis de datos nominales se enfoca exclusivamente en la frecuencia y la asociación, utilizando herramientas como las tablas de contingencia y las pruebas de independencia ([Chi-cuadrado de Pearson](#)).

El segundo subtipo fundamental es la **variable ordinal**. En este caso, las categorías sí poseen un orden natural, pero la distancia o el intervalo entre ellas no es uniforme ni conocido. Ejemplos comunes incluyen las escalas de Likert (totalmente en desacuerdo, en desacuerdo, neutro, de acuerdo, totalmente de acuerdo) o los niveles de educación (primaria, secundaria, universitaria). La naturaleza ordenada de estos datos permite el uso de métodos que aprovechan esta secuencia, como los modelos de regresión logística ordinal. Sin embargo, tratar estas variables como si fueran continuas (por ejemplo, calculando una media simple) ignora la falta de uniformidad de los intervalos y puede llevar a conclusiones erróneas. El ADC proporciona métodos específicos que respetan esta jerarquía sin asumir intervalos iguales.

Una distinción importante adicional se hace a menudo con la **variable dicotómica** o binaria, que es un caso especial de variable nominal que solo presenta dos categorías posibles (p. ej., éxito/fracaso, sí/no, vivo/muerto). Aunque es nominal, su simplicidad matemática la convierte en el foco de modelos de predicción muy potentes, como la [Regresión Logística](#). La capacidad de modelar la probabilidad de uno de dos resultados mutuamente excluyentes ha hecho de la regresión logística una de las herramientas más utilizadas en la investigación biológica y social, sirviendo como puente entre las técnicas de asociación básicas y el modelado multivariante complejo.

3. Desarrollo Histórico y Evolución Metodológica

Las raíces del Análisis de Datos Categóricos se extienden hasta principios del siglo XX, marcadas de manera indeleble por el trabajo pionero de Karl Pearson. Pearson desarrolló en 1900 la prueba de bondad de ajuste de la [Chi-cuadrado](#), que se convirtió en el pilar fundamental para evaluar si las frecuencias observadas en un conjunto de datos diferían significativamente de las frecuencias esperadas bajo la hipótesis nula de independencia o de una distribución específica. Durante las siguientes décadas, el análisis se mantuvo centrado principalmente en las pruebas de asociación bivariadas en tablas de contingencia, con estadísticos como el coeficiente de contingencia y la V

de Cramer proporcionando medidas de la fuerza de la relación. Este período inicial fue vital para establecer la necesidad de métodos no paramétricos para datos discretos.

La verdadera revolución metodológica llegó a mediados del siglo XX, impulsada por la necesidad de analizar relaciones multivariadas complejas, especialmente en encuestas sociales. La limitación del Chi-cuadrado era que solo podía manejar la asociación entre dos variables a la vez. Investigadores como Leo Goodman, Yvonne Bishop y Stephen Fienberg desarrollaron los **modelos log-lineales** a partir de la década de 1960. Estos modelos permitieron la exploración simultánea de interacciones de orden superior (tres o más variables) dentro de tablas de contingencia multidimensionales. El enfoque log-lineal transformó el ADC, permitiendo a los estadísticos modelar la estructura de las frecuencias de celdas en términos de efectos principales y efectos de interacción, de manera análoga a cómo el análisis de varianza (ANOVA) modela las medias para datos continuos.

El desarrollo culminante en el ADC moderno fue la integración de los datos categóricos dentro del marco más amplio de los **Modelos Lineales Generalizados (GLMs)**, formalizado por Nelder y Wedderburn en 1972. Los GLMs proporcionaron una estructura unificada que conecta la regresión lineal tradicional con modelos para datos no normales, como la Regresión Logística para resultados binomiales y la Regresión de Poisson para datos de conteo. Esta unificación permitió a los investigadores aplicar principios de modelado de regresión (como la inclusión de múltiples predictores, continuos o categóricos) directamente a resultados categóricos, superando las limitaciones de los modelos log-lineales puros y estableciendo la Regresión Logística como la herramienta predictiva estándar para resultados binarios y multinomiales.

4. Herramientas Analíticas Clave: Tablas de Contingencia y Pruebas de Asociación

La unidad fundamental de análisis en el ADC bivariado es la **tabla de contingencia**, también conocida como tabla cruzada. Esta matriz organiza las observaciones según las categorías de dos o más variables discretas, mostrando la frecuencia de ocurrencia de cada combinación de categorías. El examen inicial de la tabla proporciona una visión de la distribución conjunta y marginal de las variables. El objetivo principal al analizar estas tablas es determinar si las variables son independientes, es decir, si la distribución de una variable es la misma independientemente de la categoría de la otra.

Para formalizar esta evaluación de independencia, la herramienta más utilizada es la prueba de Chi-cuadrado de Pearson. Esta prueba compara las frecuencias observadas en la tabla con las frecuencias que se esperarían si las variables fueran completamente independientes. Un valor Chi-cuadrado grande, asociado con un valor p bajo, indica que la desviación de la independencia es demasiado grande para atribuirse al azar, sugiriendo una asociación significativa. Sin embargo,

es crucial recordar que la prueba Chi-cuadrado solo indica la presencia de asociación, no su dirección ni su fuerza. Además, la validez de la prueba se basa en supuestos asintóticos, requiriendo que un número suficiente de celdas tengan frecuencias esperadas mínimas (generalmente 5 o más). Cuando las frecuencias esperadas son bajas, se recurre a la **Prueba Exacta de Fisher**, que calcula la probabilidad exacta de observar la tabla dada sus marginales.

Cuando se establece la existencia de una asociación, el siguiente paso es cuantificar su fuerza y naturaleza. Para datos nominales en tablas 2x2, la medida más informativa es la **razón de probabilidades (Odds Ratio)**, especialmente en epidemiología. El Odds Ratio mide la razón de las probabilidades de ocurrencia del resultado de interés en un grupo frente a otro (por ejemplo, la probabilidad de enfermedad en expuestos frente a no expuestos). Para tablas más grandes (IxJ), existen diversas medidas de asociación, como el Tau de Kendall, la Gamma o la D de Somers, que son particularmente adecuadas para variables ordinales, ya que tienen en cuenta el orden de las categorías y miden el grado de concordancia o discordancia entre pares de observaciones.

5. Modelos de Regresión para Resultados Categóricos

El modelado de regresión constituye la faceta más avanzada del ADC, permitiendo a los investigadores examinar simultáneamente el efecto de múltiples variables predictoras (tanto categóricas como continuas) sobre un resultado categórico. La herramienta central en este ámbito es la **Regresión Logística**, diseñada específicamente para variables de respuesta binarias. A diferencia de la regresión lineal, que modela directamente la media de la variable dependiente, la regresión logística modela la transformación logarítmica de la probabilidad de éxito (el logaritmo de las probabilidades, o *logit*). Esta transformación asegura que los valores predichos de probabilidad se mantengan lógicamente dentro del rango de 0 a 1.

Cuando la variable de respuesta tiene más de dos categorías (politómica), el análisis se extiende a la **Regresión Logística Multinomial** o la **Regresión Logística Ordinal**. La regresión multinomial se aplica cuando las categorías no tienen un orden natural (nominal). Este modelo ajusta múltiples ecuaciones logísticas simultáneamente, comparando cada categoría de respuesta con una categoría de referencia seleccionada. Por otro lado, la regresión logística ordinal, también conocida como modelo de probabilidades proporcionales, se utiliza cuando las categorías tienen un orden significativo. Este modelo es más parsimonioso, ya que asume que el efecto de los predictores es constante a lo largo de las diferentes particiones de las categorías ordenadas, una suposición conocida como el supuesto de probabilidades proporcionales o pendientes paralelas.

El poder de estos modelos reside en su capacidad para estimar los coeficientes de regresión, los cuales se interpretan en términos de cambios en el logaritmo de las probabilidades (log-odds). Al exponentiarlos, estos coeficientes se convierten en **Odds Ratios ajustados**. Un Odds Ratio ajustado permite cuantificar el impacto relativo de una unidad de cambio en un predictor sobre las

probabilidades de la respuesta, manteniendo constantes los efectos de todas las demás variables en el modelo. Esta capacidad de control estadístico es fundamental para establecer relaciones causales o predictivas en estudios observacionales complejos.

6. Aplicaciones Transdisciplinarias y Relevancia Científica

La relevancia del Análisis de Datos Categóricos es inmensa y abarca virtualmente todas las disciplinas que utilizan datos de encuestas o resultados discretos. En las **Ciencias Sociales**, el ADC es indispensable para el análisis de encuestas de opinión, comportamiento electoral y actitudes sociales. Los sociólogos y politólogos utilizan la regresión logística multinomial para modelar la elección entre múltiples partidos políticos o la afiliación a diferentes grupos sociales, mientras que la regresión ordinal es clave para interpretar las respuestas a escalas de satisfacción o acuerdo. El uso de modelos log-lineales avanzados también permite desentrañar complejas interacciones entre variables demográficas y de comportamiento.

En **Epidemiología y Biometría**, el ADC es la columna vertebral de la investigación. Los estudios de casos y controles y los ensayos clínicos a menudo producen resultados binarios (p. ej., recuperación o recaída). Los Odds Ratios y los Riesgos Relativos, obtenidos a través de modelos logísticos, son las métricas estándar para cuantificar el riesgo asociado a la exposición a factores específicos. La Regresión Logística de Poisson, por otro lado, es crucial para modelar tasas de eventos (datos de conteo) como admisiones hospitalarias o número de infecciones, demostrando la versatilidad del marco GLM en la salud pública.

Finalmente, en **Investigación de Mercados y Economía**, el ADC es esencial para comprender las decisiones de los consumidores. Las empresas utilizan modelos de elección discreta (basados en la regresión multinomial) para predecir la probabilidad de que un cliente elija un producto entre varias opciones, basándose en características demográficas y atributos del producto. Esta capacidad predictiva sobre decisiones discretas tiene un impacto directo en la formulación de estrategias de fijación de precios y desarrollo de productos, consolidando el ADC como una herramienta no solo académica, sino también de alto valor empresarial.

7. Desafíos y Consideraciones Críticas

A pesar de su sofisticación, el Análisis de Datos Categóricos enfrenta varios desafíos metodológicos. Uno de los problemas más persistentes es el de la **escasez de datos** (*sparsity*), que ocurre cuando una o más celdas en la tabla de contingencia tienen frecuencias cero o muy bajas. La escasez puede invalidar los supuestos asintóticos de la prueba Chi-cuadrado y generar estimaciones de Odds Ratios infinitas o inestables en modelos logísticos. Para mitigar este problema, los estadísticos a menudo recurren a la combinación de categorías o al uso de métodos exactos, como la Regresión Logística Exacta, o técnicas de penalización bayesiana que añaden

información previa para estabilizar las estimaciones.

Otro punto de crítica importante se centra en la correcta interpretación de los resultados del modelado. Los coeficientes de la regresión logística se expresan en la escala logarítmica de las probabilidades (log-odds), lo que es inherentemente menos intuitivo que la interpretación directa de la pendiente en la regresión lineal. Si bien la transformación a Odds Ratios facilita la comunicación de los resultados, estos Odds Ratios no son equivalentes a los Riesgos Relativos, especialmente cuando la probabilidad del evento es alta. La mala interpretación de los Odds Ratios como si fueran Riesgos Relativos constituye un error común en la literatura no estadística, lo que subraya la necesidad de una cuidadosa presentación de los resultados en términos de probabilidades predichas o efectos marginales.

Finalmente, existe un debate continuo sobre la elección adecuada de métodos para variables ordinales. Aunque el uso de la regresión logística ordinal es eficiente, se basa en el estricto **supuesto de probabilidades proporcionales**. Si este supuesto es violado (es decir, el efecto de un predictor varía a través de los puntos de corte de la escala ordinal), el modelo es inapropiado y puede llevar a conclusiones sesgadas. En tales casos, los analistas deben recurrir a modelos más complejos, como el modelo de categorías adyacentes o, a menudo, tratar la variable ordinal como nominal utilizando la regresión multinomial, aunque esto conlleva la pérdida de la información sobre el orden y la necesidad de estimar un mayor número de parámetros.

8. Lecturas Adicionales

[Análisis de Datos Categóricos \(Wikipedia\)](#)

[Regresión Logística \(Wikipedia\)](#)

[Prueba Chi-cuadrado de Pearson \(Wikipedia\)](#)

[Modelo Lineal Generalizado \(Wikipedia\)](#)