

colección **textos**

# ANÁLISIS ESTADÍSTICO DE DATOS CATEGÓRICOS

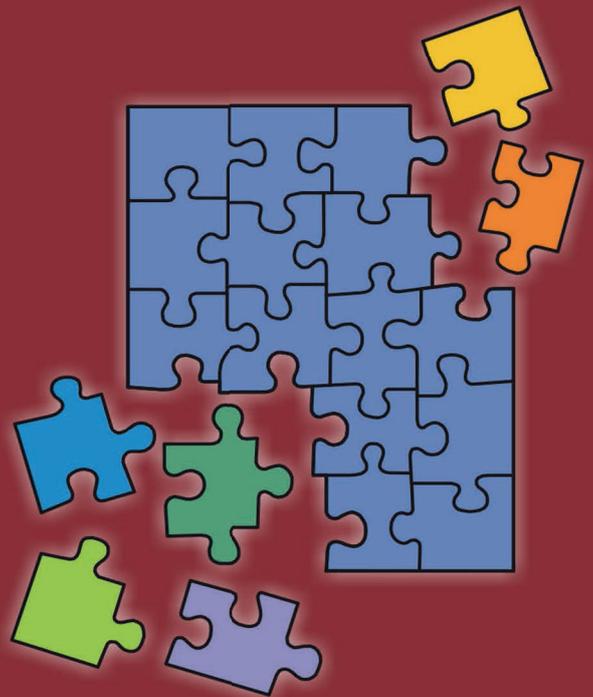
textos  
textos  
textos  
textos  
textos  
textos  
textos

Luis Guillermo Díaz Monroy  
Mario Alfonso Morales Rivera



UNIVERSIDAD NACIONAL DE COLOMBIA

SEDE BOGOTÁ  
FACULTAD DE CIENCIAS



Facultad de Ciencias  
Saber más y formar mejor



## ANÁLISIS ESTADÍSTICO DE DATOS CATEGÓRICOS



ANÁLISIS ESTADÍSTICO DE DATOS CATEGÓRICOS

LUIS GUILLERMO DÍAZ MONROY

MARIO ALFONSO MORALES RIVERA



Departamento de Estadística  
Facultad de Ciencias

Universidad Nacional de Colombia  
sede Bogotá

Análisis estadístico de datos categóricos

© Luis Guillermo Díaz Monroy  
Facultad de Ciencias  
Departamento de Estadística  
Universidad Nacional de Colombia

© Mario Alfonso Morales Rivera  
Facultad de Ciencias  
Departamento de Matemáticas y Estadística  
Universidad de Córdoba

Primera edición, 2009  
Bogotá, Colombia  
ISBN 978-958-719-186-8

Diseño de carátula: Andrea Kratzer M.

Catalogación en la publicación Universidad Nacional de Colombia

Díaz Monroy, Luis Guillermo, 1958-

Análisis estadístico de datos categóricos / Luis Guillermo Díaz Monroy,  
Mario Alfonso Morales Rivera. – Bogotá : Universidad Nacional de Co-  
lombia. Facultad de Ciencias, 2009

xvii, 376 p.

ISBN : 978-958-719-186-8

1.Análisis de regresión logística 2.Tablas de contingencia 3.Modelos log-  
lineales 4. Análisis de correspondencias (Estadística) 5. Modelos lineales  
(Estadística)

I. Morales Rivera, Mario Alfonso, 1965- II. Tít

CDD-21 519.536 / 2009

*A Daniel, Camila, Diego y Pilar, mi única categoría.*  
Luis G. Díaz

*A Nevis, mi ángel guardián.*  
Mario A. Morales



# Contenido

<b>Introducción</b>	<b>xv</b>
<b>Introducción</b>	<b>xv</b>
<b>1 Conceptos preliminares</b>	<b>1</b>
1.1 Introducción . . . . .	1
1.2 Escala de medida . . . . .	1
1.2.1 Dicotómicas . . . . .	1
1.2.2 Ordinal . . . . .	2
1.2.3 Conteos discretos . . . . .	3
1.2.4 Nominal . . . . .	3
1.3 Esquema de muestreo . . . . .	4
1.4 Modelos de muestreo . . . . .	5
1.4.1 Distribución de Poisson . . . . .	5
1.4.2 Distribución binomial . . . . .	6
1.4.3 Distribución multinomial . . . . .	8
1.4.4 Distribución hipergeométrica . . . . .	9
1.5 Inferencia sobre una proporción . . . . .	10
1.5.1 Estimación . . . . .	12
1.5.2 Distribución muestral de una proporción . . . . .	13
1.5.3 Intervalo de confianza para una proporción . . . . .	14
1.5.4 Contraste de hipótesis sobre una proporción . . . . .	15
1.6 Procesamiento de datos con R . . . . .	18

1.7	Ejercicios . . . . .	19
<b>2</b>	<b>Tablas de contingencia</b>	<b>21</b>
2.1	Introducción . . . . .	21
2.2	Tablas de contingencia . . . . .	21
2.3	Modelos probabilísticos . . . . .	26
2.3.1	Modelo de clasificación fija . . . . .	26
2.3.2	Modelo de homogeneidad . . . . .	27
2.3.3	Modelo de independencia . . . . .	29
2.4	Independencia de la clasificación . . . . .	30
2.4.1	Prueba ji-cuadrado . . . . .	31
2.4.2	Distribución ji-cuadrado . . . . .	34
2.4.3	Contraste mediante la razón de verosimilitudes ( $G^2$ )	36
2.4.4	Medidas de asociación . . . . .	37
2.4.5	Medidas ligadas a la estadística ji-cuadrado. . . . .	38
2.4.6	Medidas basadas en la reducción proporcional del error (RPE) . . . . .	39
2.4.7	Medidas de asociación ordinales . . . . .	42
2.4.8	Otras medidas de asociación . . . . .	47
2.4.9	Determinación de las fuentes de asociación . . . . .	50
2.4.10	Análisis de los residuos . . . . .	50
2.4.11	Partición de tablas . . . . .	53
2.4.12	Análisis con el PROC FREQ del paquete estadístico SAS. . . . .	55
2.5	Tablas de contingencia $2 \times 2$ . . . . .	56
2.5.1	Prueba ji-cuadrado . . . . .	56
2.5.2	La corrección por continuidad de Yates . . . . .	59
2.5.3	Prueba de la probabilidad exacta de Fisher. . . . .	59
2.5.4	Prueba de McNemar para proporciones correlacionadas en tablas $2 \times 2$ . . . . .	61
2.5.5	Riesgo relativo . . . . .	64

2.5.6	Razón de probabilidades (odds) . . . . .	67
2.5.7	Fracción etiológica . . . . .	70
2.5.8	Prueba de Cochran–Mantel–Haenszel . . . . .	72
2.6	Tablas multidimensionales . . . . .	76
2.6.1	Notación para tablas multidimensionales . . . . .	77
2.6.2	Pruebas de independencia de las variables en una ta- bla a tres vías . . . . .	78
2.6.3	Paradoja de Simpson . . . . .	79
2.7	Tamaño de muestra . . . . .	82
2.8	Procesamiento de datos con R . . . . .	85
2.9	Ejercicios . . . . .	91
<b>3</b>	<b>Análisis de correspondencias</b>	<b>93</b>
3.1	Introducción . . . . .	93
3.2	Representación geométrica de los puntos de una tabla de contingencia . . . . .	94
3.2.1	Perfiles fila y columna . . . . .	96
3.3	Semejanza entre perfiles: la distancia ji-cuadrado . . . . .	98
3.4	Explicación de la técnica . . . . .	100
3.5	Análisis de correspondencias múltiples . . . . .	104
3.5.1	Tablas de datos . . . . .	104
3.5.2	Fundamentos del análisis de correspondencias múltiples.	111
3.5.3	Propiedades del análisis de correspondencias múltiples. . . . .	112
3.5.4	Reglas de interpretación . . . . .	113
3.6	Procesamiento de datos con R . . . . .	116
3.6.1	Análisis de correspondencias simple . . . . .	116
3.6.2	Análisis de correspondencias múltiples . . . . .	119
3.7	Análisis de correspondencias múltiples mediante SAS. . . . .	120
3.8	Ejercicios . . . . .	122
<b>4</b>	<b>Modelos log–lineales</b>	<b>123</b>

4.1	Introducción . . . . .	123
4.1.1	El modelo lineal generalizado . . . . .	125
4.2	Modelos log–lineales para tablas de contingencia . . . . .	126
4.2.1	El modelo log–lineal . . . . .	126
4.2.2	Modelos jerárquicos . . . . .	129
4.2.3	Estimación de modelos log–lineales . . . . .	130
4.2.4	Ajuste de los modelos log–lineales . . . . .	133
4.2.5	Estadística ji–cuadrado de bondad de ajuste . . . . .	133
4.2.6	Residuales . . . . .	136
4.3	Procesamiento de datos con R . . . . .	141
4.4	Ejercicios . . . . .	143
<b>5</b>	<b>Regresión logística</b>	<b>145</b>
5.1	Introducción . . . . .	145
5.2	Modelo de regresión logística . . . . .	146
5.3	Interpretación de los coeficientes de regresión . . . . .	151
5.4	Construcción e interpretación de la función logística . . . . .	153
5.5	Variables ficticias (dummy) . . . . .	158
5.6	Ajuste del modelo . . . . .	162
5.6.1	Contraste de hipótesis sobre los parámetros . . . . .	162
5.6.2	Selección de modelos . . . . .	164
5.6.3	Bondad de ajuste . . . . .	169
5.7	Regresión logística con respuesta politómica . . . . .	171
5.7.1	Regresión logística nominal . . . . .	171
5.7.2	Regresión logística ordinal . . . . .	174
5.8	Algunas aplicaciones de la regresión logística . . . . .	176
5.8.1	Descripción . . . . .	177
5.8.2	Patrón de lactancia materna (LM) . . . . .	177
5.8.3	Comparación de curvas . . . . .	181

5.8.4	Índice de deserción . . . . .	185
5.8.5	Estudios prospectivos . . . . .	187
5.8.6	Estudios de cohorte . . . . .	188
5.8.7	Ensayos clínicos . . . . .	190
5.8.8	Estudios caso-control . . . . .	193
5.8.9	Razón de odds y riesgos relativos . . . . .	194
5.9	Procesamiento de datos con R . . . . .	197
5.9.1	Cálculos para la sección 5.6 . . . . .	201
5.10	Ejercicios . . . . .	204
<b>6</b>	<b>Análisis discriminante</b>	<b>207</b>
6.1	Introducción . . . . .	207
6.2	Reglas de discriminación para dos grupos . . . . .	208
6.2.1	Vía máxima verosimilitud . . . . .	208
6.3	Reglas de discriminación para varios grupos . . . . .	214
6.3.1	Grupos con matrices de covarianzas iguales . . . . .	214
6.3.2	Grupos con matrices de covarianzas distintas . . . . .	215
6.4	Tasas de error de clasificación . . . . .	217
6.4.1	Estimación de las tasas de error . . . . .	217
6.5	Otras técnicas de discriminación . . . . .	219
6.5.1	Modelo de discriminación logística para dos grupos . . . . .	219
6.5.2	Modelo de discriminación Probit . . . . .	222
6.5.3	Discriminación con datos multinomiales . . . . .	224
6.5.4	Clasificación mediante la técnica de “el vecino más cercano” . . . . .	225
6.6	Selección de variables . . . . .	226
6.7	Procesamiento de datos con R . . . . .	228
6.8	Procedimiento DISCRIM del paquete SAS . . . . .	232
6.9	Ejercicios . . . . .	233

<b>7</b>	<b>Métodos no paramétricos</b>	<b>235</b>
7.1	Introducción . . . . .	235
7.2	Pruebas de localización: una muestra . . . . .	236
7.2.1	Prueba del signo . . . . .	238
7.2.2	Muestras pareadas . . . . .	240
7.2.3	Prueba de rango signado de Wilcoxon . . . . .	242
7.3	Pruebas de localización: dos muestras . . . . .	246
7.3.1	Prueba de Mann-Whitney-Wilcoxon . . . . .	247
7.4	Pruebas de localización en diseños completamente al azar . . . . .	252
7.4.1	Prueba de Kruskal-Wallis . . . . .	253
7.5	Pruebas de localización para diseños en BAC . . . . .	256
7.5.1	Prueba de Friedman . . . . .	257
7.6	Procesamiento de datos con R . . . . .	259
7.7	Ejercicios . . . . .	263
<b>8</b>	<b>Métodos para datos de conteo</b>	<b>268</b>
8.1	Introducción . . . . .	268
8.2	Determinación de la naturaleza aleatoria de un evento . . . . .	269
8.3	Modelo de regresión tipo Poisson . . . . .	272
8.3.1	Modelo de regresión simple . . . . .	273
8.3.2	Modelo de regresión múltiple . . . . .	276
8.4	Procesamiento de datos con R . . . . .	282
8.5	Ejercicios . . . . .	284
<b>9</b>	<b>Métodos para datos emparejados</b>	<b>287</b>
9.1	Introducción . . . . .	287
9.2	Medidas de concordancia o acuerdo . . . . .	288
9.3	Estudios emparejados caso-control . . . . .	292
9.4	Regresión logística condicional . . . . .	295
9.4.1	Regresión logística simple . . . . .	296

9.4.2	Regresión logística múltiple . . . . .	299
9.5	Procesamiento de datos con R . . . . .	303
9.6	Ejercicios . . . . .	304
<b>A</b>	<b>Tablas</b>	<b>306</b>
<b>B</b>	<b>Procedimientos básicos con R</b>	<b>309</b>
B.1	Cálculo de probabilidades y cuantiles . . . . .	309
B.1.1	Distribución binomial . . . . .	310
B.1.2	Distribución de Poisson . . . . .	311
B.1.3	Distribuciones normal y ji-cuadrado . . . . .	312
B.2	Lectura de datos externos . . . . .	313
B.2.1	El directorio de trabajo . . . . .	314
B.2.2	Lectura de datos desde un archivo de texto . . . . .	314
B.2.3	Lectura de datos desde un archivo CSV . . . . .	315
B.2.4	Lectura de datos desde un archivo de Excel . . . . .	316
B.3	Selección y transformación de datos . . . . .	316
B.3.1	Creación de nuevas variables . . . . .	317
B.3.2	Selección de subconjuntos de un marco de datos . . . . .	317
B.3.3	Cálculos por niveles de un factor . . . . .	320
	<b>Bibliografía</b>	<b>322</b>
	<b>Índice temático</b>	<b>328</b>

# Tablas

1.1	Resultados respiratorios. . . . .	2
1.2	Datos de artritis. . . . .	2
1.3	Niños con problemas respiratorios. . . . .	3
1.4	Tipo de sangre por región de procedencia. . . . .	3
1.5	Distribución de Poisson y binomial con $\mu = 2.0$ . . . . .	8
2.1	Opinión sobre el servicio de salud. . . . .	21
2.2	Tabla de contingencia. . . . .	22
2.3	Tabla de contingencia completa (de la tabla 2.1). . . . .	24
2.4	Evaluación de un funcionario. . . . .	27
2.5	Concepto sobre el aborto. . . . .	28
2.6	Drogas vs. prácticas bisexuales. . . . .	29
2.7	Frecuencias esperadas. . . . .	35
2.8	Opinión sobre el servicio de salud. . . . .	40
2.9	Opinión sobre el servicio de salud (de tabla 2.1). . . . .	44
2.10	Concordancias. . . . .	45
2.11	Discordancias. . . . .	45
2.12	Residuales. . . . .	52
2.13	Salida SAS. . . . .	57
2.14	Tabla de contingencia $2 \times 2$ . . . . .	58
2.15	Resultados respiratorios. . . . .	58
2.16	Curación de infecciones severas. . . . .	59
2.17	Probabilidades de las tablas $2 \times 2$ . . . . .	61

2.18	Frecuencias de muestras apareadas. . . . .	62
2.19	Recuperación en pacientes depresivos. . . . .	63
2.20	Sujetos que muestran náusea con las drogas $A$ y $B$ . . . . .	64
2.21	Consumo de aspirina e infartos del miocardio. . . . .	66
2.22	Consumo de aspirina e infartos del miocardio. . . . .	71
2.23	Mejoría en enfermedades respiratorias. . . . .	75
2.24	Tabla de contingencia tridimensional. . . . .	77
2.25	Enfermedades cardiacas por tabaquismo y edad. . . . .	79
2.26	Edad entre 25 y 45 años. . . . .	80
2.27	Edad superior a 45 años. . . . .	80
2.28	Tabaquismo y enfermedades cardiacas. . . . .	80
2.29	Admisiones a una universidad por género. . . . .	81
2.30	Admisiones por facultad y género. . . . .	81
2.31	Datos sobre accidentes automovilísticos. . . . .	91
2.32	Datos sobre uso de marihuana por estudiantes. . . . .	91
2.33	Comparación entre radiación y cirugía en el tratamiento de cáncer de laringe. . . . .	92
2.34	Datos de dolor tras la cirugía. . . . .	92
3.1	Frecuencias absolutas. . . . .	94
3.2	Frecuencia relativas. . . . .	94
3.3	Perfil fila. . . . .	97
3.4	Perfil columna. . . . .	97
3.5	Color de ojos vs. color del cabello. . . . .	101
3.6	Coordenadas, color de ojos y del cabello. . . . .	102
3.7	Coordenadas y contribuciones de las modalidades. . . . .	115
3.8	Respuesta de la enfermedad de Hodgkin a un tratamiento según la tipología. . . . .	122
4.1	Datos de melanoma maligno. . . . .	131
4.2	Valores esperados ( $E_{ij}^*$ ) de los datos de la tabla 4.1. . . . .	131
4.3	Parámetros estimados (PROC CATMOD). . . . .	132

4.4	Tabla de análisis de varianza (PROC CATMOD).	135
4.5	Parámetros estimados para el modelo (4.18).	136
4.6	Datos sobre enfermedades coronarias.	138
4.7	Parámetros estimados para el modelo (4.19).	138
4.8	Parámetros estimados para el modelo (4.20).	139
4.9	Parámetros estimados para el modelo (4.21).	140
4.10	Raza y pena de muerte.	144
4.11	Úlceras gástrica y duodenal en relación con el uso de aspirina.	144
5.1	Infección en pacientes hospitalizados.	147
5.2	Pacientes por modelo de atención y condición de infección.	147
5.3	Enfermedades coronarias frente a tabaquismo, edad y TAS.	155
5.4	Estimaciones máximo verosímiles con los datos de la tabla 5.3.	158
5.5	Pacientes por grupo sanguíneo, RH y condición patológica.	160
5.6	Verificación de los parámetros, $H_0 : \beta_i = 0$ .	164
5.7	Summary of Stepwise Procedure	167
5.8	Summary of Backward Elimination Procedure	167
5.9	Datos de artritis.	173
5.10	Verificación de los parámetros, $H_0 : \beta_i = 0$ .	176
5.11	Valores de la función logística.	180
5.12	Estimación según régimen de atención primaria.	181
5.13	Deserción de lactancia materna en los primeros tres meses para cuatro subpoblaciones.	185
5.14	Cohorte de 2.000 pacientes infartados.	188
5.15	Decesos por tabaquismo.	188
5.16	Decesos por edad.	189
5.17	Modelos ajustados para 2.000 infartados.	190
5.18	Esquema de datos sobre un ensayo clínico de acupuntura.	192
5.19	Ajuste de la probabilidad de mejoría.	193

5.20	Resultados de un estudio caso-control para evaluar letalidad en infartados con hábito de fumar y edad como factores explicativos. . . . .	195
5.21	Modelos ajustados para 400 casos y 400 controles. . . . .	196
5.22	Pacientes que se mueven o quejan al hacer una incisión 15 minutos después de aplicada la concentración del anestésico. . . . .	204
5.23	Datos de inhibición. . . . .	205
6.1	Evaluación psiquiátrica. . . . .	211
6.2	Datos de acupuntura. . . . .	233
7.1	Distribución de $T^+$ con $n = 4$ . . . . .	244
7.2	Distribución de $U$ con $n_1 = 3$ y $n_2 = 2$ . . . . .	248
7.3	Hipótesis alternativas y regiones de rechazo, prueba de Mann-Whitney. . . . .	249
7.4	Consumo de cloruro de sodio. . . . .	252
7.5	Datos sobre variación de pesos de pacientes tratados para várices. . . . .	264
7.6	Tiempos para desarrollar una tarea con o sin alcohol. . . . .	265
7.7	Niveles de NDMA. . . . .	265
7.8	Niveles de alquitrán. . . . .	266
7.9	Niveles de plaguicida (ppb). . . . .	266
7.10	Reducción de peso en libras. . . . .	267
8.1	Frecuencias observadas y esperadas. . . . .	271
8.2	Casos nuevos de melanomas. . . . .	275
8.3	Regresión ajustada a los casos con melanomas. . . . .	276
8.4	Datos sobre cáncer en la piel. . . . .	280
8.5	Estimación del modelo de regresión múltiple con los datos de la tabla 8.4. . . . .	281
8.6	Razón de verosimilitud. . . . .	282
8.7	Número de pólizas de seguros y número de reclamos. . . . .	285
8.8	Muertes por enfermedades coronarias. . . . .	286

9.1	Concordancia entre dos observadores. . . . .	289
9.2	Probabilidades de concordancia entre dos observadores . . .	289
9.3	Diagnóstico de dos neurólogos. . . . .	291
9.4	Proporciones factor $\times$ enfermedad. . . . .	292
9.5	Frecuencias caso $\times$ control. . . . .	293
9.6	Diagnóstico previo de diabetes para MI. Pares caso-control.	295
9.7	Emparejamiento $1 : m_i$ . . . . .	296
9.8	Bajo peso al nacer. . . . .	298
9.9	Emparejamiento $n_i : m_i$ . . . . .	299
9.10	Estimación para datos peso bajo al nacer. . . . .	302
9.11	Influencia de los anticonceptivos orales sobre el cáncer endo- metrial. . . . .	305
A.1	Distribución normal acumulada . . . . .	307
A.2	Percentiles de la distribución ji-cuadrado. . . . .	308

# Figuras

1.1	Distribución binomial. . . . .	7
1.2	Esquematización de una distribución hipergeométrica. . . . .	10
1.3	Función de verosimilitud para $X = 0, 4$ y $8$ -éxitos. . . . .	12
1.4	Región de rechazo para $H_0 : \pi = p_0$ . . . . .	17
2.1	Distribución de frecuencias. . . . .	25
2.2	Perfiles fila de la opinión por estrato. . . . .	26
2.3	Región de rechazo para $H_0 : \pi = p_0$ . . . . .	33
2.4	Valores de la razón de odds (RO). . . . .	69
3.1	Tabla de frecuencias y sus marginales. . . . .	96
3.2	Perfiles fila. . . . .	98
3.3	Perfiles columna. . . . .	99
3.4	Representación de los datos de color de ojos ( $\Delta$ ) y cabello ( $\times$ )	103
3.5	Esquema del análisis de correspondencias . . . . .	105
3.6	Tabla múltiple. . . . .	107
3.7	Construcción de la tabla de Burt. . . . .	109
3.8	Variables activas y suplementarias en el plano factorial . . . . .	117
5.1	Función logística . . . . .	149
5.2	Curva de prevalencia de lactancia materna. . . . .	179
5.3	Curvas de prevalencia de lactancia materna por modelos de atención. . . . .	182
5.4	Curvas de prevalencia de lactancia materna. Ajuste bivariado	183

5.5	Curvas de prevalencia de consumo de cuatro alimentos. . .	186
6.1	Discriminación lineal. . . . .	211
6.2	Discriminación en senil o no senil. . . . .	212
6.3	Regiones de discriminación para tres grupos. . . . .	216
6.4	Función logística. . . . .	220
6.5	Discriminación probit. . . . .	223
7.1	Distribución sesgada de mediana 0. . . . .	243

# Introducción

La distinción entre los llamados datos *cualitativos* y los denominados *cuantitativos* no siempre es clara, pues en algunos casos variables de tipo cuantitativo pueden considerarse como variables categóricas al dividir su rango de valores en intervalos o categorías, esto corresponde a una categorización de una variable cuantitativa. Un tratamiento recíproco puede considerarse para las variables cualitativas, es decir que pueden transformarse a variables cuantitativas, este procedimiento se muestra con el análisis de correspondencias. En estas notas se hace una revisión, bastante panorámica, sobre algunas metodologías estadísticas que coadyuvan al esclarecimiento e interpretación de la información contenida en datos categóricos.

El texto ha sido elaborado pensando en un lector que demande el uso de algunas herramientas estadísticas, útiles para el análisis de la información, principalmente de tipo categórico, producto de algún trabajo de investigación. No obstante que los primeros destinatarios son las personas que trabajen en torno a problemas de la salud y la biología, el material estadístico que se ofrece puede ser empleado por investigadores de otras disciplinas, pues basta cambiar el escenario de los ejemplos e ilustraciones, para hacer de este texto un instrumento de apoyo a varias disciplinas.

La primera parte contiene algunos conceptos generales junto con el tratamiento clásico de datos categóricos a través del análisis de tablas de contingencia, los cuales se desarrollan en los capítulos 1 y 2. Tratamientos alternativos a las tablas de contingencia se desarrollan en los capítulos 3 y 4 mediante el análisis de correspondencias y los modelos log-lineales. En el capítulo 5 se presenta el modelamiento con variable respuesta de tipo categórico el cual se hace a través de la regresión logística. El capítulo 6 contiene algunas de las técnicas de discriminación de uso más frecuente. En el capítulo 7 se esquematizan algunos contrastes de tipo no paramétrico sobre estadísticas de localización. Se tratan, en el capítulo 8, algunas técnicas estadísticas para datos de conteo. Finalmente, en el capítulo 9, se desarrolla la técnica de emparejamiento de datos.

Para el desarrollo de los cálculos que las estimaciones y estadísticas requie-

ren, se hace uso, principalmente, de los paquetes SAS y R. En cada capítulo se presenta la sintaxis pertinente para la ejecución de tales cómputos. Se debe advertir que existen otras herramientas computacionales igualmente útiles, tales como SPSS, BMDP, MINITAB, STATA, entre otras.

Agradecemos al Grupo de Investigación en “Estadística aplicada en la investigación experimental, la industria y la biotecnología”, al Departamento de Estadística de la Universidad Nacional y al Departamento de Matemáticas y Estadística de la Universidad de Córdoba por posibilitar y permitir el ofrecimiento de estas notas.

Las notas se deben principalmente a la bibliografía que se anexa al final y a las preguntas, comentarios y sugerencias de nuestros colegas y de nuestros estudiantes. Este es un material que se puede mejorar en la medida que sea leído y cuestionado, por tanto agradecemos los comentarios y sugerencias que surjan de su estudio.

*Luis Guillermo Díaz Monroy*

*Mario Alfonso Morales Rivera*

# Capítulo 1

## Conceptos preliminares

### 1.1 Introducción

En este capítulo se presentan los aspectos fundamentales, a manera de elementos estadísticos básicos, para el desarrollo de los demás temas. Se revisa la naturaleza de los datos categóricos, el modelo probabilístico *binomial*, el de *Poisson* y la inferencia sobre una *proporción*.

### 1.2 Escala de medida

La escala de medida de una variable categórica es un elemento importante para la selección del análisis estadístico apropiado. Una selección inadecuada de la escala de medida puede conducir a una estrategia estadística inapropiada que arrojaría conclusiones erróneas acerca de la realidad contenida en los datos.

#### 1.2.1 Dicotómicas

Son variables que tienen dos posibles respuestas, frecuentemente corresponden a la presencia o no de cierto atributo. Por ejemplo: ¿Desarrolló el sujeto la enfermedad? ¿En los últimos tres meses ha fumado alguna vez, o no? ¿Está afiliado actualmente al régimen contributivo de salud, o no?, etc.

Por ejemplo, el objetivo de un ensayo clínico para un nuevo medicamento contra la gripe es saber si los pacientes alivian sus dolencias. La tabla 1.1

contiene información sobre 124 pacientes, quienes recibieron tratamiento (medicamento) o un placebo (sin medicamento).

Tabla 1.1: Resultados respiratorios.

Tratamiento	Favorable	Desfavorable	Total
Placebo	16	48	64
Prueba	40	20	60

El grupo placebo consta de 64 pacientes; mientras que el grupo de prueba del medicamento contiene 60 pacientes.

### 1.2.2 Ordinal

En muchas ocasiones las variables categóricas representan más de dos posibles resultados, y a veces estos resultados poseen un orden propio. Tales variables tienen una escala de medida *ordinal*.

El estado de mejoría o progreso de un paciente se puede calificar como marcado (3), regular (2), ninguno (1). Este es el caso de un ensayo clínico en el que se investiga un tratamiento para la artritis reumatoidea. A hombres y mujeres les fue asignada una actividad (tratamiento) o un placebo (no actividad). Se midió el nivel de progreso o mejoría conseguido al final del ensayo; los datos están dispuestos en la tabla 1.2.

Tabla 1.2: Datos de artritis.

Sexo	Tratamiento	Progreso			Total
		Marcado	Regular	Ninguno	
Femenino	Actividad	16	5	6	27
Femenino	Placebo	6	7	19	32
Masculino	Actividad	5	2	7	14
Masculino	Placebo	1	0	10	11

*Fuente:* Stokes, Davis y Koch (1997: 218)

Note que las variables categóricas pueden manejarse de diferentes formas. Por ejemplo, en la tabla 1.2 se pueden fusionar las columnas Marcado y Regular para producir una variable dicotómica: “Progreso” frente a “No progreso”. Este tipo de agrupamiento se hace generalmente durante el análisis cuando hay interés por esta clase de respuestas o cuando se quiere obtener información adicional sobre los datos.

### 1.2.3 Conteos discretos

Corresponde a los casos en los que en lugar de registrar categorías, los resultados son números enteros. Por ejemplo, una investigación sobre enfermedades respiratorias en niños de diferentes zonas visitados dos veces determina si ellos mostraron síntomas de la enfermedad. La respuesta medida fue si los niños exhibieron síntomas en 0, 1 o 2 periodos. La tabla 1.3 contiene estos resultados.

Tabla 1.3: Niños con problemas respiratorios.

Sexo	Tratamiento	Periodos			Total
		0	1	2	
Femenino	Rural	45	64	71	180
Femenino	Urbana	80	104	116	300
Masculino	Rural	84	124	82	290
Masculino	Urbana	106	117	87	310

### 1.2.4 Nominal

Si se dispone de variables con más de dos categorías, a las cuales no se les atribuye un orden, se tiene una variable de tipo *nominal*. Por ejemplo, el tipo de sangre de una persona y la región geográfica donde nació; la tabla 1.4 muestra esta información. En este tipo de variables la relación

Tabla 1.4: Tipo de sangre por región de procedencia.

Región	Tipo de sangre				Total
	O	A	B	AB	
Norte	40	160	150	50	400
Centro	50	130	100	40	320
Sur	70	180	90	60	400
Total	160	470	340	150	<b>1.120</b>

que se puede establecer entre sus categorías es estrictamente de igualdad (o desigualdad).

### 1.3 Esquema de muestreo

Cuando el interés en el estudio es de tipo inferencial, los datos categóricos pueden proceder de diferentes esquemas de muestreo, sea este probabilístico o no. La naturaleza del muestreo determina los supuestos que pueden hacerse para desarrollar y aplicar un análisis estadístico determinado. Generalmente, los datos se ubican en uno de tres esquemas muestrales: datos históricos, datos experimentales y datos de encuestas.

Los datos históricos se refieren a estudios en los cuales los datos tienen una definición geográfica o circunstancial. Por ejemplo, la ocurrencia de una enfermedad infecciosa en una área determinada, los niños atendidos en un centro de salud, o el número de accidentes durante un periodo específico.

Los datos experimentales son extraídos de estudios que involucran la asignación aleatoria de tratamientos a un grupo de sujetos. Es el caso en el que a los sujetos se les administra una dosis entre varias dosis de un medicamento. En ciencias de la salud, los datos experimentales pueden incluir pacientes a quienes se les administra un placebo o un medicamento de acuerdo con sus condiciones médicas.

En estudios por encuestas, los individuos son seleccionados aleatoriamente desde una población objetivo. Por ejemplo, se selecciona una muestra de los usuarios de determinado medicamento para investigar algunos rasgos físicos de estos. El investigador puede seleccionar aleatoriamente una población de estudio y luego asignar aleatoriamente tratamientos a los individuos que resulten para el estudio.

La principal diferencia entre los tres esquemas de muestreo asociados a los ejemplos anteriores es el empleo de la aleatorización para obtenerlos. Los datos históricos no involucran aleatorización; en consecuencia, es difícil asumir que ellos representan determinada población. Los datos experimentales tienen una buena cobertura de la población, la cual está restringida por los tratamientos considerados en el protocolo del estudio; en el muestreo por encuestas, los datos tienen una muy buena cobertura de alguna población grande.

La unidad de aleatorización (en conexión con la unidad de observación) puede ser un sujeto o un conglomerado de sujetos. Además, la aleatorización puede aplicarse a sujetos, llamados estratos o bloques, con igual o desigual probabilidad. En muestreo por encuestas, esto puede conducir a diseños complejos, como el muestreo aleatorio estratificado, o un diseño de conglomerados por múltiples etapas. En estudios de diseño experimental, tales consideraciones llevan a estudios de medidas repetidas, datos longitudinales, entre otros.

## 1.4 Modelos de muestreo

El análisis de datos categóricos, o casi cualquier tipo de análisis estadístico, requiere supuestos acerca del mecanismo de aleatorización que genera los datos; esto es, el modelo probabilístico desde el cual se asume que son generados los datos. Se presentan a continuación las distribuciones de probabilidad de uso más frecuente en el análisis de datos categóricos.

### 1.4.1 Distribución de Poisson

La distribución de Poisson es un modelo probabilístico adecuado para evaluar la probabilidad de ocurrencia de un evento en intervalos de tiempo, longitud, superficie o volumen. Por ejemplo, el número de accidentes por semana en un tramo de carretera, el número de aves muertas por kilómetro cuadrado en una región. Si  $X$  es la variable aleatoria que cuenta el número de veces que un evento ocurre por intervalo (tiempo, longitud, superficie, etc), y si se tiene que el número promedio de eventos por intervalo es  $\mu$ , las probabilidades de los posibles resultados  $k = 0, 1, 2, \dots$  se calculan mediante la expresión

$$P(X = k) = \frac{e^{-\mu} \mu^k}{k!}, \quad \text{para } k = 0, 1, 2, \dots \quad (1.1)$$

El término  $k!$  es llamado el *factorial* de  $k$  y denota el producto de los  $k$ -primeros enteros, es decir,  $k! = 1 \times 2 \times 3 \times \dots \times k$ , con  $0! = 1$ . El término  $e^{-\mu}$  denota la *función exponencial*, algunas veces expresada como  $\exp(-\mu)$ ; siendo  $e \approx 2.7182$  el cual es la base de los *logaritmos naturales*. Esto último significa que  $e^a = b$  si y solo si  $\ln(b) = a$ . Por ejemplo  $e^{0.7} = \exp(0.7) = 2.0$  corresponde a que  $\ln(2.0) = 0.7$ .

Suponga que el número de personas con infarto que acuden a una clínica tiene una tasa promedio de 2 por semana. Mediante el modelo de Poisson con  $\mu = 2$ , (i) la probabilidad de 0 infartos ( $k = 0$ ), y, (ii) de a lo más un infarto, en una semana cualquiera, por (1.1), es igual, respectivamente, a:

$$\begin{aligned} (i) \quad P(X = 0) &= \frac{e^{-2} 2^0}{0!} = e^{-2} = 0.1353 \\ (ii) \quad P(X \leq 1) &= \frac{e^{-2} 2^0}{0!} + \frac{e^{-2} 2^1}{1!} = 0.1353 + 0.2707 = 0.4060. \end{aligned}$$

donde  $P(X = 1) = e^{-2}(2^1)/1! = 0.2707$  es la probabilidad que se presente un infarto durante una semana. De (i) la probabilidad de que hayan infartos es:  $P(X \geq 1) = 1 - P(X = 0) = 1 - 0.1353 = 0.8647$

El valor esperado de la variable aleatoria  $X$  con distribución de Poisson es igual a su varianza; es decir,

$$E(X) = \text{var}(X) = \mu, \quad \sigma(X) = \sqrt{\mu} \quad (1.2)$$

Para el caso de los infartos por semana, si la tasa de ocurrencia de estos permanece constante de una semana a otra, entonces en un periodo largo el conteo de estos tendría una media de alrededor de 2 y una desviación estándar cercana a  $\sqrt{2} = 1.41$ .

De acuerdo con los parámetros dados en (1.2) se observa que la varianza se incrementa a medida que la media lo hace; los conteos tienden a variar más cuando el nivel de sus promedios es alto. Así, cuando el número de infartos por semana es 10, se observa que la variabilidad es más grande que cuando el número es igual a 2 por semana.

## 1.4.2 Distribución binomial

En el ejemplo anterior, el número de infartos fatales semanal es aleatorio. El número de infartos semanal, fatales o no, es también aleatorio. En muchas aplicaciones se tiene como fijo el número de veces que se presenta un fenómeno. En cada caso, el resultado es un evento  $A$  o no es el evento  $A$  (es  $A^c$ ); entonces, se quiere registrar las veces que este fenómeno ocurre con la característica determinada ( $A$ ) en un número fijo de observaciones ( $n$ ).

Un ejemplo es el caso de infarto fatal ( $A$ ) o no fatal ( $A^c$ ) y, en general, la ocurrencia o no de un evento; también es común hablar de “éxito” o “fracaso” para referirse al evento  $A$  o al  $A^c$ , respectivamente.

Una variable aleatoria  $X$  tiene distribución *binomial* si su función de probabilidad está dada por

$$P(X = k) = \binom{n}{k} \pi^k (1 - \pi)^{n-k}, \text{ para } k = 0, 1, 2, \dots, n. \quad (1.3)$$

donde los dos parámetros  $n$  y  $\pi = P(A)$  son tales que  $n$  es un entero no negativo y  $0 \leq \pi \leq 1$ . Se escribe  $X \sim B(n, \pi)$ . Para  $n = 1$  la variable aleatoria se denomina de *Bernoulli*; es decir, que una variable aleatoria binomial es una suma de variables independientes tipo Bernoulli.

La cantidad  $\binom{n}{k} = \frac{n!}{k!(n-k)!}$  es el número de posibles arreglos de  $k$ -elementos (subconjuntos) que se pueden formar a partir de un conjunto que tiene  $n$ -elementos. Así, por ejemplo, el número de formas como se puede conformar un comité de 3 personas, escogidas de un grupo de 5 personas, es

$$\binom{5}{3} = \frac{5!}{3!2!} = 10$$