

**UDS**

**ANTOLOGIA**

UDS

Antología

# Estadística Inferencial

Licenciatura en Administración y Estrategias de Negocios

Licenciatura en Nutrición

4o Cuatrimestre



## Marco Estratégico de Referencia

### Antecedentes Históricos

Nuestra Universidad tiene sus antecedentes de formación en el año de 1979 con el inicio de actividades de la normal de educadoras “Edgar Robledo Santiago”, que en su momento marcó un nuevo rumbo para la educación de Comitán y del estado de Chiapas. Nuestra escuela fue fundada por el Profesor de Primaria Manuel Albores Salazar con la idea de traer Educación a Comitán, ya que esto representaba una forma de apoyar a muchas familias de la región para que siguieran estudiando.

En el año 1984 inicia actividades el CBTiS Moctezuma Ilhuicamina, que fue el primer bachillerato tecnológico particular del estado de Chiapas, manteniendo con esto la visión en grande de traer Educación a nuestro municipio, esta institución fue creada para que la gente que trabajaba por la mañana tuviera la opción de estudiar por las tardes.

La Maestra Martha Ruth Alcázar Mellanes es la madre de los tres integrantes de la familia Albores Alcázar que se fueron integrando poco a poco a la escuela formada por su padre, el Profesor Manuel Albores Salazar; Víctor Manuel Albores Alcázar en septiembre de 1996 como chofer de transporte escolar, Karla Fabiola Albores Alcázar se integró como Profesora en 1998, Martha Patricia Albores Alcázar en el departamento de finanzas en 1999.

En el año 2002, Víctor Manuel Albores Alcázar formó el Grupo Educativo Albores Alcázar S.C. para darle un nuevo rumbo y sentido empresarial al negocio familiar y en el año 2004 funda la Universidad Del Sureste.

La formación de nuestra Universidad se da principalmente porque en Comitán y en toda la región no existía una verdadera oferta Educativa, por lo que se veía urgente la creación de una institución de Educación superior, pero que estuviera a la altura de las exigencias

de los jóvenes que tenían intención de seguir estudiando o de los profesionistas para seguir preparándose a través de estudios de posgrado.

Nuestra Universidad inició sus actividades el 18 de agosto del 2004 en las instalaciones de la 4a avenida oriente sur no. 24, con la licenciatura en Puericultura, contando con dos grupos de cuarenta alumnos cada uno. En el año 2005 nos trasladamos a nuestras propias instalaciones en la carretera Comitán – Tzitol km. 57 donde actualmente se encuentra el campus Comitán y el Corporativo UDS, este último, es el encargado de estandarizar y controlar todos los procesos operativos y Educativos de los diferentes Campus, Sedes y Centros de Enlace Educativo, así como de crear los diferentes planes estratégicos de expansión de la marca a nivel nacional e internacional.

Nuestra Universidad inició sus actividades el 18 de agosto del 2004 en las instalaciones de la 4a avenida oriente sur no. 24, con la licenciatura en Puericultura, contando con dos grupos de cuarenta alumnos cada uno. En el año 2005 nos trasladamos a nuestras propias instalaciones en la carretera Comitán – Tzitol km. 57 donde actualmente se encuentra el campus Comitán y el corporativo UDS, este último, es el encargado de estandarizar y controlar todos los procesos operativos y educativos de los diferentes campus, así como de crear los diferentes planes estratégicos de expansión de la marca.

### Misión

Satisfacer la necesidad de Educación que promueva el espíritu emprendedor, aplicando altos estándares de calidad Académica, que propicien el desarrollo de nuestros alumnos, Profesores, colaboradores y la sociedad, a través de la incorporación de tecnologías en el proceso de enseñanza-aprendizaje.

### Visión

Ser la mejor oferta académica en cada región de influencia, y a través de nuestra Plataforma Virtual tener una cobertura Global, con un crecimiento sostenible y las ofertas académicas innovadoras con pertinencia para la sociedad.

### Valores

- Disciplina
- Honestidad
- Equidad
- Libertad

---

## Escudo



El escudo de la UDS, está constituido por tres líneas curvas que nacen de izquierda a derecha formando los escalones al éxito. En la parte superior está situado un cuadro motivo de la abstracción de la forma de un libro abierto.

## Eslogan

“Mi Universidad”

## Albores



Es nuestra mascota, un Jaguar. Su piel es negra y se distingue por ser líder, trabaja en equipo y obtiene lo que desea. El ímpetu, extremo valor y fortaleza son los rasgos que distinguen.





## Materia

Objetivo de la materia:

Que el alumno consolide la competencia habilitante de la lectura y escritura al reconocer y ejercer las cuatro habilidades de la lengua: escuchar, leer, hablar y escribir, con el fin de aplicarlas a diversas situaciones de su vida, académicas y cotidianas.

Criterios y procedimientos de evaluación y acreditación:

Actividad en plataforma	30%
Tareas	10%
Examen	60%
Total	100%
Escala de calificaciones	6 - 10
Mínima aprobatoria	7





## Índice general

1	Introducción a la Estadística Inferencial .....	17
1.1	Breve historia de la estadística .....	17
1.2	Concepto de estadística .....	19
1.3	Estadística descriptiva .....	19
1.3.1	Medidas de tendencia central .....	20
1.3.2	Medidas de dispersión .....	20
1.3.3	Distribución de frecuencias .....	21
1.3.4	Gráficos y representaciones visuales .....	21
1.3.5	Análisis de datos en nutrición .....	21
1.3.6	Importancia de la estadística descriptiva en nutrición .....	22
1.4	Estadística inferencial .....	22
1.4.1	Muestras y poblaciones .....	22
1.4.2	Estimación .....	23
1.4.3	Contraste de hipótesis .....	23
1.4.4	Importancia de la estadística inferencial en nutrición .....	23
1.4.5	Errores en la inferencia estadística .....	24
1.4.6	Aplicación de la estadística inferencial en nutrición .....	24
1.5	Introducción a la inferencia estadística .....	24
1.5.1	Diferencias entre estadística descriptiva e inferencial .....	25
1.5.2	Componentes clave de la inferencia estadística .....	25
1.5.3	Procedimientos en inferencia estadística .....	26
1.5.4	Importancia de la inferencia estadística en nutrición .....	26

---

1.6	Teoría de decisiones en estadística	27
1.6.1	Elementos fundamentales de la teoría de decisiones	28
1.6.2	Criterios de decisión	28
1.6.3	Árboles de decisión	29
1.6.4	Aplicaciones de la teoría de decisiones en nutrición	29
1.7	Componentes de una investigación estadística	30
1.7.1	Planteamiento del problema	30
1.7.2	Revisión de la literatura	30
1.7.3	Diseño de la investigación	30
1.7.4	Población y muestra	31
1.7.5	Recolección de datos	31
1.7.6	Análisis de datos	31
1.7.7	Interpretación de resultados	31
1.7.8	Comunicación de resultados	31
1.8	Recolección de datos	32
1.8.1	Fuentes de datos	32
1.8.2	Métodos de recolección de datos	32
1.8.3	Consideraciones éticas	33
1.8.4	Validación de datos	33
1.9	Estadística paramétrica	33
1.9.1	Fundamentos de la estadística paramétrica	34
1.9.2	Condiciones de aplicabilidad	34
1.9.3	Ejemplos de métodos estadísticos paramétricos	34
1.9.4	Ventajas y desventajas	35
1.10	Población	35
1.10.1	Definición de población	36
1.10.2	Características de la población	36
1.10.3	Tipos de poblaciones	36
1.10.4	Importancia de la población en la investigación nutricional	36
1.10.5	Ejemplos de poblaciones en nutrición	37
1.11	Muestra aleatoria	37
1.11.1	Definición de muestra aleatoria	37
1.11.2	Importancia de la muestra aleatoria	37
1.11.3	Tipos de muestreo aleatorio	38
1.11.4	Métodos de selección de muestra aleatoria	38
1.11.5	Ejemplos de muestreo aleatorio en nutrición	38
2	Inferencia Estadística: Estimación, Muestreo	41
2.1	Teoría de conjuntos	41
2.1.1	Conjuntos	41
2.1.2	Operaciones con conjuntos	41
2.1.3	Diagramas de Venn	41

---

2.1.4	Definición de Conjunto . . . . .	42
2.1.5	Tipos de Conjuntos . . . . .	42
2.1.6	Notación de Conjuntos . . . . .	42
2.1.7	Operaciones con Conjuntos . . . . .	43
2.1.8	Propiedades de los Conjuntos . . . . .	43
2.1.9	Aplicaciones en Estadística . . . . .	43
2.2	Distribución de Muestreo . . . . .	43
2.2.1	Definición . . . . .	44
2.2.2	Teorema del Límite Central . . . . .	44
2.2.3	Importancia de la Distribución de Muestreo . . . . .	44
2.2.4	Ejemplo de Distribución de Muestreo . . . . .	45
2.2.5	Distribuciones de Muestreo Comunes . . . . .	45
2.2.6	Consideraciones sobre la Muestra . . . . .	45
2.3	Muestreo Aleatorio Simple . . . . .	45
2.3.1	Definición . . . . .	46
2.3.2	Proceso de Muestreo . . . . .	46
2.3.3	Ejemplo en Nutrición . . . . .	46
2.3.4	Ventajas del Muestreo Aleatorio Simple . . . . .	46
2.3.5	Desventajas del Muestreo Aleatorio Simple . . . . .	46
2.4	Muestreo Aleatorio Estratificado Simple . . . . .	47
2.4.1	Definición . . . . .	47
2.4.2	Características del Muestreo Aleatorio Estratificado Simple . . . . .	47
2.4.3	Pasos para implementar el Muestreo Aleatorio Estratificado Simple . . . . .	47
2.4.4	Ventajas del Muestreo Aleatorio Estratificado Simple . . . . .	48
2.4.5	Desventajas . . . . .	48
2.4.6	Consideraciones . . . . .	48
2.4.7	Ejemplo de Muestreo Aleatorio Estratificado Simple . . . . .	48
2.4.8	Aplicaciones típicas . . . . .	49
2.5	Muestreo por Conglomerado . . . . .	50
2.5.1	Definición . . . . .	50
2.5.2	Características del Muestreo por Conglomerado . . . . .	50
2.5.3	Tipos de Muestreo por Conglomerado . . . . .	50
2.5.4	Ventajas del Muestreo por Conglomerado . . . . .	50
2.5.5	Desventajas del Muestreo por Conglomerado . . . . .	51
2.5.6	Pasos para Implementar el Muestreo por Conglomerado . . . . .	51
2.5.7	Ejemplo de Muestreo por Conglomerado . . . . .	51
2.5.8	Aplicaciones Comunes del Muestreo por Conglomerado . . . . .	52
2.5.9	Fórmula del Error Estándar en el Muestreo por Conglomerado . . . . .	52
2.5.10	Consideraciones Finales . . . . .	52
2.6	Intervalo de Confianza para la Diferencia entre Medias . . . . .	52
2.6.1	Definición . . . . .	52
2.6.2	Fórmula . . . . .	53
2.6.3	Ejemplo de Aplicación . . . . .	53

---

2.7	Muestreo Estratificado	54
2.7.1	Definición	54
2.7.2	Tipos de Muestreo Estratificado	54
2.7.3	Ventajas del Muestreo Estratificado	54
2.7.4	Desventajas del Muestreo Estratificado	55
2.7.5	Pasos para Implementar el Muestreo Estratificado	55
2.7.6	Fórmula para el Tamaño de Muestra en Muestreo Estratificado Proporcional	55
2.7.7	Ejemplo de Muestreo Estratificado	55
2.8	Principio Aditivo, Multiplicativo y Arreglo Rectangular	55
2.8.1	Definición	55
2.8.2	Principio Aditivo	55
2.8.3	Principio Multiplicativo	56
2.8.4	Arreglo Rectangular	56
2.8.5	Relación entre Principio Aditivo, Multiplicativo y Arreglo Rectangular	57
2.8.6	Ejemplo Combinado de Principio Aditivo y Multiplicativo	57
2.9	Diagrama de Árbol y Principio Multiplicativo	57
2.9.1	Definición del Principio Multiplicativo	57
2.9.2	Definición del Diagrama de Árbol	57
2.9.3	Construcción de un Diagrama de Árbol	58
2.9.4	Aplicaciones del Diagrama de Árbol	58
2.9.5	Relación entre Diagrama de Árbol y Principio Multiplicativo	59
2.9.6	Ejemplo de Aplicación Combinada	59
2.9.7	Ventajas de Usar Diagramas de Árbol	59
2.9.8	Desventajas de Usar Diagramas de Árbol	59
2.9.9	Aplicaciones Prácticas del Principio Multiplicativo y los Diagramas de Árbol	59
2.10	Permutaciones	60
2.10.1	Definición	60
2.10.2	Fórmula para Permutaciones	60
2.10.3	Ejemplo de Permutaciones	60
2.10.4	Permutaciones de un Conjunto Completo	61
2.10.5	Permutaciones con Elementos Repetidos	61
2.10.6	Aplicaciones de las Permutaciones	61
2.10.7	Relación entre Permutaciones y Combinaciones	62
2.11	Combinaciones	62
2.11.1	Definición	62
2.11.2	Fórmula para Combinaciones	62
2.11.3	Ejemplo de Combinaciones	62
2.11.4	Combinaciones de un Conjunto Completo	63
2.11.5	Combinaciones con Elementos Repetidos	63
2.11.6	Aplicaciones de las Combinaciones	63
2.11.7	Relación entre Combinaciones y Permutaciones	64
2.11.8	Ejemplo de Aplicación Combinada	64

---

3	Asociación Estadística entre Variables	65
3.1	Asociación Estadística entre Variables	65
3.1.1	¿Qué es una asociación entre variables?	65
3.1.2	Diferencia entre asociación y causalidad	65
3.1.3	Conceptos clave	66
3.2	Midiendo la Asociación entre Dos Variables	66
3.2.1	Coefficiente de correlación de Pearson	66
3.2.2	Correlación de Spearman	67
3.2.3	Interpretación de las correlaciones	68
3.3	El caso de dos variables categóricas	68
3.3.1	¿Qué son las variables categóricas?	68
3.3.2	Tablas de contingencia	68
3.3.3	Medidas de asociación para dos variables categóricas	69
3.3.4	Interpretación	69
3.4	El caso de una variable categórica y una continua	69
3.4.1	¿Qué son las variables continuas?	69
3.4.2	Comparación de medias	70
3.4.3	Prueba t de Student	70
3.4.4	Interpretación	71
3.4.5	Conclusión	71
3.5	El caso de dos variables cuantitativas	71
3.5.1	¿Qué son las variables cuantitativas?	71
3.5.2	Diagrama de dispersión	72
3.5.3	Coefficiente de correlación	72
3.6	El modelo de regresión lineal	73
3.6.1	¿Qué es la regresión lineal?	73
3.6.2	Estimación de los parámetros: método de los mínimos cuadrados	73
3.6.3	Interpretación del modelo	74
3.6.4	Evaluación del modelo: coeficiente de determinación $R^2$	74
3.6.5	Concepto de error o residual	75
3.6.6	Propiedades importantes del análisis de regresión lineal	75
3.7	Bondad de ajuste del modelo de regresión	75
3.7.1	¿Qué es la bondad de ajuste?	75
3.7.2	Error estándar de los residuos	76
3.7.3	Conclusión sobre la bondad de ajuste	76
3.8	Teoría de la probabilidad	76
3.8.1	Definiciones Básicas	77
3.8.2	Probabilidad Clásica	77
3.8.3	Propiedades de la Probabilidad	77

---

3.9	Modelos teóricos de distribución de probabilidad	77
3.9.1	Distribuciones Discretas	77
3.9.2	Distribuciones Continuas	77
3.10	La distribución binomial	77
3.10.1	Función de probabilidad	78
3.11	La distribución o curva normal	78
3.11.1	Función de densidad de probabilidad	78
3.11.2	Ejemplo	78
3.12	La selección de la muestra	79
3.12.1	¿Qué es la selección de la muestra?	79
3.12.2	Tipos de muestreo	79
3.12.3	Tamaño de la muestra	79
4	Prueba de Hipótesis	81
4.1	Metodología para la prueba de hipótesis	81
4.2	Hipótesis nula y alternativa	82
4.2.1	Hipótesis nula ( $H_0$ )	82
4.2.2	Hipótesis alternativa ( $H_1$ )	83
4.3	Error tipo I y tipo II	83
4.3.1	Error de tipo I	84
4.3.2	Error de tipo II	84
4.3.3	Ejemplo	84
4.4	Pruebas de hipótesis	84
4.5	Prueba de hipótesis Z para la media	84
4.5.1	Requisitos para realizar la Prueba Z	85
4.5.2	Cálculo de la Prueba Z y Ejemplo con Distribuciones Normales	85
4.5.3	Tabla de Valores Críticos Z para la Prueba de Hipótesis	85
4.5.4	Tabla de Valores Críticos Z para la Prueba de Dos Colas	86
4.5.5	Tabla de Valores Críticos Z para la Prueba de Una Cola	86
4.5.6	Caso 1: Prueba de dos colas	86
4.5.7	Caso 2: Prueba de una cola (a la derecha)	87
4.5.8	Ejemplo 1: Prueba de Hipótesis para la Media (Prueba Z Unilateral)	87
4.5.9	Prueba Z y Valor Crítico	87
4.5.10	Ejemplo 2: Prueba de Hipótesis para la Media (Prueba Z Bilateral)	87
4.6	Varianza	88
4.6.1	Utilidad de la varianza	89
4.6.2	Varianza Poblacional	89
4.6.3	Varianza Muestral	90

---

	15	
4.7	Desviación Estándar	90
4.7.1	Utilidad de la desviación estándar . . . . .	91
4.7.2	Desviación Estándar Poblacional . . . . .	92
4.7.3	Desviación Estándar Muestral . . . . .	92
4.8	Ejemplo detallado de la Prueba Z para la Media Poblacional	93
	Bibliografía . . . . .	96





## 1 — Introducción a la Estadística Inferencial

### 1.1 Breve historia de la estadística

La estadística tiene raíces antiguas, comenzando con registros de censos y datos económicos en civilizaciones como las de Egipto y Babilonia. En el campo de la nutrición, el uso de la estadística comenzó a cobrar relevancia en el siglo XX, con la aparición de estudios sobre la relación entre alimentación y salud. Investigadores como Ancel Keys, a través del Estudio de los Siete Países, utilizaron métodos estadísticos para entender cómo los hábitos alimentarios afectaban la prevalencia de enfermedades como la cardiopatía coronaria.

La estadística, como disciplina, ha tenido una evolución larga y diversa, desarrollándose en diferentes contextos a lo largo de los siglos. Su origen puede rastrearse a antiguas civilizaciones que recogían datos numéricos sobre la población, las tierras y la riqueza, pero su formalización como ciencia moderna llegó mucho más tarde. Aquí presentamos un recorrido por los principales hitos históricos de la estadística y su relevancia actual, especialmente en el campo de la nutrición.

#### Orígenes en la Antigüedad

Los primeros registros de recolección de datos numéricos provienen de civilizaciones como las de Babilonia, Egipto y China, que utilizaban estos datos para censos y fines administrativos. Por ejemplo, los egipcios realizaron censos para organizar la mano de obra en la construcción de monumentos, y los chinos recogían datos para administrar el impuesto agrícola. En la Antigua Grecia, Heródoto describió el uso de estadísticas para determinar la población y la riqueza de las naciones que visitaba.

#### La estadística en la Edad Media

Durante la Edad Media, la estadística siguió usándose principalmente para la administración de recursos y censos de población. Sin embargo, no fue hasta el Renacimiento cuando se vio

un renacer del pensamiento científico y la cuantificación de fenómenos. En el siglo XV, los registros detallados de nacimientos, defunciones y matrimonios comenzaron a usarse con mayor rigor en ciudades europeas, lo que sentó las bases para el análisis demográfico.

#### Siglo XVII: Primeros desarrollos formales

En el siglo XVII, la estadística comenzó a tomar forma como ciencia. El trabajo de John Graunt, quien analizó los registros de mortalidad en Londres, es uno de los primeros ejemplos de lo que hoy llamamos estadísticas descriptivas. Graunt es considerado el padre de la demografía moderna, ya que fue el primero en utilizar datos numéricos para hacer observaciones sobre las causas de la muerte y la esperanza de vida.

Otro avance crucial fue la teoría de la probabilidad, desarrollada por matemáticos como Blaise Pascal y Pierre de Fermat, que sentó las bases para el análisis estadístico en situaciones de incertidumbre. Esta teoría sería fundamental más adelante para el desarrollo de la estadística inferencial.

#### Siglo XVIII: Consolidación de la estadística

Durante el siglo XVIII, la estadística empezó a aplicarse en campos como la medicina y la agricultura. Se comenzaron a desarrollar métodos más precisos para la recolección y análisis de datos, que facilitaron la planificación y evaluación de políticas públicas. En esta época, el término “estadística” empezó a utilizarse formalmente, derivado del latín *status*, refiriéndose al “estado” o “gobierno”, ya que en sus inicios esta disciplina se centraba en datos relacionados con la administración pública.

#### Siglo XIX: La era de la estadística moderna

El siglo XIX fue un período de grandes avances en la estadística moderna, con la creación de instituciones y el desarrollo de técnicas estadísticas más avanzadas. Karl Pearson y Francis Galton fueron figuras clave en este proceso, aportando conceptos como la correlación y la regresión, que permiten estudiar la relación entre variables.

En el campo de la nutrición, la estadística se utilizó para estudiar la relación entre la dieta y las enfermedades. Un ejemplo notable es el trabajo de William Farr, quien analizó la relación entre la ingesta alimentaria y las tasas de mortalidad, poniendo de relieve la importancia de una buena nutrición para la salud pública.

#### Siglo XX: Estadística aplicada en nutrición

Con el auge de la estadística inferencial y el desarrollo de la computación en el siglo XX, la estadística se convirtió en una herramienta esencial en campos como la epidemiología, la biomedicina y la nutrición. Investigadores como Ancel Keys utilizaron análisis estadísticos para demostrar la relación entre la dieta mediterránea y una menor incidencia de enfermedades cardiovasculares, lo que sentó las bases para la investigación nutricional moderna.

En la actualidad, la estadística es fundamental en la evaluación de la seguridad alimentaria,

los estudios de intervenciones dietéticas y la planificación de políticas públicas relacionadas con la salud nutricional.

#### Relevancia actual

Hoy en día, la estadística sigue siendo un pilar en el campo de la nutrición. Desde la evaluación de hábitos alimentarios hasta la determinación de las necesidades nutricionales de poblaciones específicas, los métodos estadísticos permiten a los profesionales de la salud y nutrición tomar decisiones basadas en evidencia. El análisis estadístico de grandes bases de datos de salud, como encuestas nacionales de consumo alimentario, es vital para identificar tendencias y áreas de mejora en la nutrición pública.

El desarrollo de la estadística ha estado estrechamente ligado a la evolución de la ciencia y la tecnología. En nutrición, la aplicación de métodos estadísticos ha permitido avanzar en el conocimiento sobre cómo la alimentación afecta la salud y el bienestar. A medida que la ciencia continúa desarrollándose, la estadística seguirá siendo una herramienta indispensable para responder a nuevas preguntas y desafíos en el campo de la nutrición.

## 1.2 Concepto de estadística

La estadística es una rama de la matemática aplicada que se ocupa de la recolección, análisis, interpretación y presentación de datos. En nutrición, se emplea para analizar patrones dietéticos, evaluar intervenciones nutricionales y estudiar la relación entre la ingesta de alimentos y los resultados de salud. La estadística nos ayuda a tomar decisiones basadas en evidencia científica.

## 1.3 Estadística descriptiva

La estadística descriptiva se enfoca en resumir los datos de manera que se puedan entender fácilmente. En nutrición, esto podría incluir la presentación de datos sobre el consumo de calorías diarias promedio en una población, o la distribución de macronutrientes en diferentes grupos de estudio. Herramientas como tablas, gráficos de barras, histogramas y medidas como la media, la mediana y la desviación estándar permiten simplificar los datos para su análisis.

Por ejemplo, en un estudio sobre el consumo de frutas y verduras en adultos, podríamos usar la media para describir el número promedio de porciones consumidas diariamente, y la desviación estándar para mostrar la variabilidad entre los individuos.

La estadística descriptiva es la rama de la estadística que se encarga de recolectar, organizar, presentar y resumir un conjunto de datos, permitiendo describir sus características principales de manera clara y comprensible. En el campo de la nutrición, la estadística descriptiva es fundamental para entender patrones alimentarios, características de grupos poblacionales y la relación entre la dieta y la salud. A continuación, se presentan los conceptos clave y las herramientas utilizadas en la estadística descriptiva.

### 1.3.1 Medidas de tendencia central

Las medidas de tendencia central son valores que representan el centro o punto medio de un conjunto de datos. Las más comunes son la media, la mediana y la moda.

Media

La media aritmética, también conocida como promedio, es el valor obtenido al sumar todos los datos y dividir el resultado entre el número de observaciones. Se utiliza frecuentemente para describir características poblacionales, como el consumo promedio de calorías diarias en un grupo de individuos.

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad (1.1)$$

Donde  $\bar{x}$  es la media,  $x_i$  son los valores observados, y  $n$  es el número total de observaciones.

Mediana

La mediana es el valor que divide a un conjunto de datos ordenados en dos partes iguales. En estudios nutricionales, la mediana es útil cuando los datos presentan valores extremos, ya que no se ve afectada por estos, a diferencia de la media.

Moda

La moda es el valor que aparece con mayor frecuencia en un conjunto de datos. En nutrición, la moda puede utilizarse para identificar el alimento más consumido por un grupo de personas.

### 1.3.2 Medidas de dispersión

Las medidas de dispersión nos indican cuán dispersos o extendidos están los datos en torno a su tendencia central. Las medidas más comunes son el rango, la varianza y la desviación estándar.

Rango

El rango es la diferencia entre el valor máximo y el valor mínimo de un conjunto de datos. En nutrición, puede utilizarse para identificar la amplitud en el consumo de un determinado nutriente en una población.

$$Rango = x_{\text{máx}} - x_{\text{mín}} \quad (1.2)$$

Varianza y desviación estándar

La varianza mide la dispersión de los datos en torno a la media, mientras que la desviación estándar es la raíz cuadrada de la varianza, lo que permite interpretar la dispersión en las mismas unidades que los datos originales.

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n} \quad (1.3)$$

$$\sigma = \sqrt{\sigma^2} \quad (1.4)$$

Donde  $\sigma^2$  es la varianza,  $\sigma$  es la desviación estándar,  $x_i$  son los valores observados, y  $\bar{x}$  es la media.

En estudios de nutrición, la desviación estándar nos indica cuán variable es el consumo de calorías o nutrientes en una población. Por ejemplo, una desviación estándar alta en el consumo de sodio podría sugerir que algunas personas consumen significativamente más sodio que otras.

### 1.3.3 Distribución de frecuencias

La distribución de frecuencias organiza los datos en categorías o intervalos y muestra la frecuencia con que cada uno de ellos aparece. Es una herramienta fundamental en nutrición para agrupar datos como el número de personas que consumen una determinada cantidad de frutas o vegetales.

La distribución de frecuencias puede representarse gráficamente utilizando histogramas, polígonos de frecuencia, o diagramas de barras. Estas representaciones permiten visualizar la distribución de los datos de manera clara.

### 1.3.4 Gráficos y representaciones visuales

La representación visual de los datos es una de las herramientas más poderosas de la estadística descriptiva. En nutrición, se utilizan gráficos para comunicar información de manera eficaz y comprensible. Entre los gráficos más comunes se encuentran:

- Gráficos de barras: Utilizados para comparar categorías de datos. Por ejemplo, se pueden emplear para mostrar el consumo de diferentes grupos de alimentos (cereales, lácteos, vegetales) en una población.
- Histogramas: Son útiles para representar la distribución de una variable cuantitativa. Un histograma del consumo de calorías diarias podría mostrar cómo se distribuye la población en función de su ingesta calórica.
- Diagramas de pastel (o torta): Se emplean para representar proporciones. Por ejemplo, un diagrama de pastel podría ilustrar qué porcentaje del consumo total de energía proviene de carbohidratos, proteínas y grasas en la dieta de un grupo de personas.

### 1.3.5 Análisis de datos en nutrición

En estudios de nutrición, la estadística descriptiva es clave para identificar patrones dietéticos y diferencias entre subgrupos de población. A través de la organización y resumen de los datos, los investigadores pueden observar tendencias generales, como la prevalencia de deficiencias nutricionales o el exceso de consumo de ciertos nutrientes.

Por ejemplo, el análisis de la distribución del índice de masa corporal (IMC) en una población puede ayudar a identificar los niveles de obesidad o desnutrición. Asimismo, las medidas descriptivas pueden revelar cómo varía el consumo de micronutrientes esenciales, como el hierro o el calcio, entre diferentes grupos de edad o sexo.

### 1.3.6 Importancia de la estadística descriptiva en nutrición

La estadística descriptiva es el primer paso en el análisis de datos y proporciona una visión clara de los patrones y características de los datos antes de aplicar técnicas más avanzadas. En nutrición, es crucial para resumir grandes volúmenes de datos de encuestas alimentarias, ensayos clínicos, y estudios poblacionales. Estos resúmenes permiten identificar problemas clave de salud pública y diseñar intervenciones basadas en datos sólidos.

Por ejemplo, el análisis de las encuestas nacionales de salud y nutrición (ENSANUT) en México ha permitido a las autoridades identificar problemas de desnutrición infantil, así como patrones de sobrepeso y obesidad en adultos, proporcionando una base de datos fundamental para la creación de políticas públicas.

## 1.4 Estadística inferencial

La estadística inferencial se utiliza para hacer generalizaciones o predicciones sobre una población, basadas en los datos recolectados de una muestra. En el campo de la nutrición, se emplea para evaluar si una intervención, como una dieta baja en grasas, tiene un efecto significativo sobre la reducción del colesterol en una población.

A través de métodos inferenciales, como pruebas de hipótesis y estimaciones, se puede determinar si las diferencias observadas en los datos son estadísticamente significativas o si podrían haber ocurrido por azar.

La estadística inferencial es la rama de la estadística que permite hacer generalizaciones o predicciones sobre una población a partir de la información obtenida de una muestra. A diferencia de la estadística descriptiva, que se enfoca en resumir y describir datos específicos, la estadística inferencial busca hacer afirmaciones más amplias y responder preguntas que van más allá de los datos observados.

En el contexto de la nutrición, la estadística inferencial es crucial para evaluar el estado nutricional de poblaciones, determinar la efectividad de intervenciones dietéticas y hacer recomendaciones de salud pública. A continuación, se detallan los conceptos clave de la estadística inferencial y su importancia en el análisis de datos nutricionales.

### 1.4.1 Muestras y poblaciones

Uno de los conceptos centrales de la estadística inferencial es la distinción entre población y muestra. La población incluye a todos los individuos o elementos que comparten una característica particular que queremos estudiar, como el total de niños en un país que sufren de desnutrición. Sin embargo, en muchos casos es imposible o impráctico recopilar datos de todos los individuos de una población. Aquí es donde entra en juego la muestra, que es un subconjunto representativo de la población.

En nutrición, la muestra puede consistir en un grupo de personas seleccionadas para participar en un estudio sobre hábitos alimentarios. A partir de esta muestra, los nutricionistas

pueden inferir conclusiones sobre los hábitos alimentarios de toda la población.

#### 1.4.2 Estimación

La estimación es uno de los principales objetivos de la estadística inferencial. A partir de los datos de la muestra, los investigadores hacen inferencias sobre parámetros desconocidos de la población, como la media del consumo de calorías diarias o la prevalencia de deficiencias nutricionales.

Existen dos tipos principales de estimación: la estimación puntual y la estimación por intervalo. En la estimación puntual, se proporciona un único valor como mejor estimación del parámetro poblacional. En la estimación por intervalo, se ofrece un rango de valores, lo que proporciona una idea de la incertidumbre en la estimación.

Por ejemplo, en un estudio sobre el consumo de frutas y verduras, se puede utilizar la estadística inferencial para estimar cuántas porciones de estos alimentos consume, en promedio, una población a partir de los datos de una muestra representativa.

#### 1.4.3 Contraste de hipótesis

El contraste de hipótesis es otra herramienta fundamental de la estadística inferencial. Este proceso permite a los investigadores evaluar afirmaciones sobre la población basadas en los datos de una muestra. En nutrición, esto podría implicar probar si una dieta específica tiene un efecto significativo en la pérdida de peso o si la suplementación con vitaminas reduce la prevalencia de enfermedades.

El contraste de hipótesis comienza con una afirmación o hipótesis sobre un parámetro poblacional, llamada hipótesis nula. A continuación, se recogen datos de una muestra y se utilizan métodos estadísticos para determinar si la evidencia es suficiente para rechazar la hipótesis nula en favor de una hipótesis alternativa.

Por ejemplo, un investigador puede querer comprobar si una nueva dieta baja en carbohidratos reduce significativamente el colesterol en comparación con una dieta tradicional. Aquí, la hipótesis nula sería que no hay diferencia entre las dos dietas, mientras que la hipótesis alternativa sugeriría que la dieta baja en carbohidratos tiene un efecto mayor.

#### 1.4.4 Importancia de la estadística inferencial en nutrición

La estadística inferencial es fundamental para avanzar en el conocimiento científico en el campo de la nutrición. Permite a los investigadores tomar decisiones basadas en datos limitados, pero representativos, y aplicar los resultados a grupos más amplios de personas. En nutrición, esto es clave para la toma de decisiones sobre recomendaciones dietéticas, políticas de salud pública y estrategias de intervención para mejorar el estado nutricional.

Por ejemplo, los estudios nutricionales que utilizan estadística inferencial pueden ayudar a determinar la cantidad óptima de nutrientes necesarios para prevenir deficiencias en dife-

rentes grupos de población. Además, la estadística inferencial puede utilizarse para evaluar el impacto de programas de suplementación o fortificación alimentaria en la reducción de enfermedades nutricionales, como la anemia o la deficiencia de vitamina A.

#### 1.4.5 Errores en la inferencia estadística

Dado que la estadística inferencial se basa en muestras, siempre existe la posibilidad de cometer errores en las conclusiones. Hay dos tipos principales de errores en la inferencia estadística:

- Error de tipo I: Ocurre cuando se rechaza la hipótesis nula cuando en realidad es verdadera. En nutrición, esto podría implicar concluir que una dieta es efectiva cuando en realidad no lo es.
- Error de tipo II: Ocurre cuando no se rechaza la hipótesis nula cuando en realidad es falsa. En nutrición, esto podría significar no detectar el efecto beneficioso de un suplemento nutricional cuando realmente existe.

Los nutricionistas deben tener en cuenta estos posibles errores al interpretar los resultados de estudios basados en muestras, y utilizar procedimientos adecuados para minimizar estos riesgos.

#### 1.4.6 Aplicación de la estadística inferencial en nutrición

En la práctica, la estadística inferencial permite a los profesionales de la salud y nutrición tomar decisiones basadas en datos empíricos. Algunas aplicaciones comunes incluyen:

- Estudios clínicos: En los estudios clínicos sobre intervenciones dietéticas, la estadística inferencial se utiliza para determinar si una intervención (como un cambio en la dieta o la suplementación con un nutriente específico) tiene un efecto significativo sobre los resultados de salud.
- Encuestas poblacionales: En encuestas nacionales de salud y nutrición, la estadística inferencial permite hacer inferencias sobre el estado nutricional de toda una población basándose en los datos de una muestra representativa.
- Evaluación de riesgos: En estudios de seguridad alimentaria, la estadística inferencial es fundamental para evaluar el riesgo de deficiencias o excesos de nutrientes en poblaciones, y para planificar programas de intervención eficaces.

La estadística inferencial es una herramienta esencial en el análisis de datos, especialmente en campos como la nutrición, donde no siempre es posible analizar datos de toda la población. Gracias a las técnicas inferenciales, los investigadores pueden tomar decisiones informadas y hacer recomendaciones con un alto grado de confianza, mejorando la salud y el bienestar de las poblaciones a través de la dieta y la nutrición.

### 1.5 Introducción a la inferencia estadística

La inferencia estadística nos permite extraer conclusiones sobre una población basándonos en una muestra de datos. En nutrición, esto es fundamental para realizar estudios a gran escala que proporcionen información sobre la dieta y la salud. Por ejemplo, en un estudio de intervención nutricional, podemos tomar una muestra representativa de individuos, medir su ingesta calórica y luego inferir si una dieta específica es efectiva para reducir el peso en

la población general.

La inferencia estadística es un conjunto de técnicas que permite a los investigadores sacar conclusiones sobre una población basándose en los datos obtenidos de una muestra. Su objetivo principal es proporcionar herramientas para generalizar las observaciones de una muestra a toda la población y tomar decisiones en situaciones de incertidumbre. En otras palabras, permite realizar predicciones y formular hipótesis sobre una población más amplia a partir de la información limitada que proporciona una muestra.

La inferencia estadística es fundamental en estudios de nutrición, ya que permite evaluar el estado nutricional de grandes poblaciones sin tener que examinar a cada individuo. A través de técnicas inferenciales, los nutricionistas y profesionales de la salud pueden hacer recomendaciones dietéticas basadas en muestras representativas de personas y aplicar estos resultados a la población en general.

### 1.5.1 Diferencias entre estadística descriptiva e inferencial

La estadística descriptiva se enfoca en resumir y describir los datos recolectados de manera directa, utilizando medidas como la media, la mediana, la moda, la varianza y los percentiles. Sin embargo, no permite hacer conclusiones más allá de los datos observados.

Por el contrario, la estadística inferencial permite extender las conclusiones obtenidas de una muestra a la población de interés, ayudando a contestar preguntas más generales. Esto es particularmente útil cuando trabajar con toda la población es impráctico o imposible, como en estudios nutricionales a gran escala donde no es posible analizar a cada persona de una región o país.

### 1.5.2 Componentes clave de la inferencia estadística

La inferencia estadística se basa en varios conceptos fundamentales que permiten hacer generalizaciones fiables:

- **Muestra y población:** La población es el conjunto total de individuos o elementos que se desea estudiar, mientras que la muestra es un subconjunto representativo de esa población. Por ejemplo, en estudios sobre el consumo de frutas y verduras, la población puede ser un país entero, mientras que la muestra puede estar compuesta por unas pocas miles de personas.
- **Estimación:** Consiste en utilizar los datos de una muestra para calcular valores aproximados de los parámetros de la población, como la media o la proporción. En el ámbito nutricional, esto puede implicar estimar la ingesta media de calorías diarias en una población a partir de los datos obtenidos de una muestra de individuos.
- **Contraste de hipótesis:** Permite evaluar si una afirmación específica sobre la población es consistente con los datos observados. Esto es útil en nutrición para probar si una dieta o intervención tiene efectos significativos sobre la salud. Por ejemplo, se podría contrastar la hipótesis de que una dieta baja en carbohidratos reduce el colesterol en comparación con una dieta convencional.

- Error muestral: Dado que la inferencia estadística se basa en muestras y no en la población completa, siempre existe un grado de incertidumbre. El error muestral es la diferencia entre un parámetro de la muestra (como la media) y el parámetro correspondiente en la población.

### 1.5.3 Procedimientos en inferencia estadística

Existen dos métodos principales que se emplean en la inferencia estadística: la estimación y el contraste de hipótesis. A continuación, se describen brevemente ambos enfoques y su relevancia en el campo de la nutrición.

#### Estimación

La estimación en inferencia estadística se refiere al proceso de inferir el valor de un parámetro poblacional a partir de los datos de la muestra. Hay dos tipos de estimaciones comunes:

- Estimación puntual: Proporciona un único valor como mejor estimación del parámetro poblacional. Por ejemplo, en un estudio sobre hábitos alimenticios, la estimación puntual puede ser la media de consumo diario de calorías en una muestra de individuos.
- Estimación por intervalo: Proporciona un rango de valores, llamado intervalo de confianza, que tiene una alta probabilidad de contener el parámetro poblacional. Este método es útil cuando se quiere tener en cuenta la incertidumbre en la estimación. En un estudio de nutrición, un intervalo de confianza puede proporcionar una mejor idea de la ingesta media de ciertos nutrientes en una población.

#### Contraste de hipótesis

El contraste de hipótesis es una técnica de inferencia estadística que se utiliza para tomar decisiones o realizar afirmaciones sobre una población basándose en los datos de una muestra. El proceso consiste en plantear dos hipótesis:

- Hipótesis nula ( $H_0$ ): Es una afirmación inicial que supone que no hay efecto o diferencia. Por ejemplo, en un estudio sobre el impacto de una dieta baja en sodio en la presión arterial, la hipótesis nula podría ser que la dieta no tiene ningún efecto en la presión arterial.
- Hipótesis alternativa ( $H_a$ ): Es la afirmación opuesta que sugiere que sí existe un efecto o diferencia. En el mismo ejemplo, la hipótesis alternativa podría ser que la dieta baja en sodio reduce la presión arterial.

El resultado del contraste de hipótesis permite a los investigadores decidir si los datos son lo suficientemente fuertes como para rechazar la hipótesis nula en favor de la hipótesis alternativa. En el campo de la nutrición, esta técnica es utilizada frecuentemente para evaluar la efectividad de distintas dietas o suplementos nutricionales.

### 1.5.4 Importancia de la inferencia estadística en nutrición

La inferencia estadística es vital en estudios de nutrición porque permite a los investigadores extraer conclusiones sobre el comportamiento alimentario, las deficiencias nutricionales o el impacto de dietas a partir de muestras relativamente pequeñas. Esto es crucial en contextos

donde obtener datos de toda una población sería extremadamente costoso o imposible.

Por ejemplo, los estudios epidemiológicos que evalúan la prevalencia de la obesidad o la deficiencia de vitaminas en una población suelen basarse en muestras representativas. Con la ayuda de la inferencia estadística, los resultados obtenidos de estas muestras se pueden generalizar a toda la población, lo que facilita la implementación de políticas de salud pública y el desarrollo de guías alimentarias.

La inferencia estadística es una herramienta poderosa que permite a los investigadores hacer generalizaciones sobre una población a partir de una muestra, lo que es esencial en estudios de nutrición y salud pública. A través de la estimación y el contraste de hipótesis, los profesionales de la nutrición pueden tomar decisiones basadas en datos, mejorando así la planificación y ejecución de intervenciones dietéticas y políticas de salud nutricional. A medida que el campo de la nutrición sigue creciendo, la inferencia estadística continuará siendo un recurso clave para avanzar en la comprensión de los efectos de los hábitos alimentarios en la salud.

## 1.6 Teoría de decisiones en estadística

La teoría de decisiones en estadística nos ayuda a tomar decisiones informadas cuando hay incertidumbre. En nutrición, esto podría aplicarse en situaciones en las que se debe decidir si una intervención alimentaria tiene un impacto positivo en la salud de un grupo determinado. Mediante la asignación de probabilidades y el análisis de riesgos, los nutricionistas pueden tomar decisiones sobre recomendaciones dietéticas basadas en datos.

Por ejemplo, si un estudio sugiere que reducir el consumo de sodio puede disminuir la presión arterial, la teoría de decisiones nos permite evaluar si esta reducción debería recomendarse a toda la población o solo a grupos específicos.

La teoría de decisiones en estadística es un campo que se centra en la toma de decisiones bajo incertidumbre. Se utiliza para evaluar y elegir entre diferentes acciones posibles basadas en los resultados observados y la información disponible, a menudo en situaciones donde no es posible tener certeza completa sobre los resultados futuros. En el contexto de la estadística, la teoría de decisiones proporciona un marco para tomar decisiones basadas en datos, teniendo en cuenta tanto la variabilidad como el riesgo inherente en los procesos de muestreo y estimación.

En estudios de nutrición, la teoría de decisiones puede ser aplicada para determinar el mejor curso de acción en situaciones como la elección de tratamientos dietéticos, la recomendación de suplementos nutricionales o la identificación de intervenciones en salud pública. Este enfoque es crucial para tomar decisiones informadas que maximicen los beneficios y minimicen los riesgos.

### 1.6.1 Elementos fundamentales de la teoría de decisiones

La teoría de decisiones se basa en varios componentes clave que guían el proceso de toma de decisiones. Estos elementos son aplicables en una variedad de contextos, incluidos los estudios de nutrición, y se describen a continuación:

- **Decisiones:** El primer paso es identificar las decisiones que deben ser tomadas. En nutrición, esto podría ser, por ejemplo, decidir si implementar o no un programa de educación alimentaria en una comunidad.
- **Estados de la naturaleza:** Se refiere a las posibles condiciones o circunstancias que afectan el resultado de una decisión, pero sobre las que el tomador de decisiones no tiene control. Estos estados pueden ser conocidos o desconocidos. En nutrición, un estado de la naturaleza podría ser el nivel de deficiencia de ciertos nutrientes en una población que aún no ha sido evaluado completamente.
- **Acciones:** Las acciones son los cursos de acción que un decisor puede tomar. En el caso de un nutricionista, las acciones podrían incluir la recomendación de una dieta específica, la administración de suplementos, o la promoción de hábitos alimenticios saludables.
- **Consecuencias:** Cada decisión tiene una consecuencia asociada, que puede depender de las condiciones bajo las cuales se tomó la decisión. Por ejemplo, si se recomienda una dieta baja en calorías para una población con alto riesgo de obesidad, las consecuencias pueden incluir la reducción del peso corporal, pero también efectos secundarios como una deficiencia de energía si no se gestiona adecuadamente.
- **Probabilidades:** A menudo, las consecuencias de una decisión están vinculadas a eventos inciertos. En muchos casos, la probabilidad de que un determinado estado de la naturaleza ocurra no es completamente conocida. En nutrición, por ejemplo, se podría desconocer exactamente cuántos individuos de una población tienen una deficiencia de vitamina D, pero se pueden estimar probabilidades basadas en estudios previos.
- **Utilidad:** Se refiere a la medición subjetiva de la satisfacción o valor que un tomador de decisiones asocia con una consecuencia particular. En nutrición, la utilidad puede estar relacionada con la mejora del bienestar general de una población, la reducción de enfermedades relacionadas con la alimentación o el aumento de la calidad de vida.

### 1.6.2 Criterios de decisión

En la teoría de decisiones existen varios criterios que pueden utilizarse para seleccionar la mejor acción en función de las probabilidades y consecuencias esperadas. Estos criterios son importantes cuando se necesita tomar decisiones informadas en situaciones de incertidumbre, como es el caso de muchos estudios de nutrición.

- **Criterio de maximización de la utilidad esperada:** Este criterio sugiere que se debe elegir la acción que maximice la utilidad esperada, es decir, el valor esperado de las consecuencias ponderadas por las probabilidades de los diferentes estados de la naturaleza. En estudios nutricionales, podría implicar recomendar la dieta que tenga el mayor beneficio promedio en una población, considerando tanto los beneficios nutricionales como los posibles riesgos.
- **Criterio de minimización del riesgo:** En algunos casos, es preferible minimizar el ries-

go en lugar de maximizar el beneficio esperado. Este enfoque es útil en nutrición cuando se trabaja con poblaciones vulnerables o en condiciones donde los riesgos de intervenciones equivocadas son altos. Por ejemplo, si existe incertidumbre sobre los efectos a largo plazo de una intervención dietética, se podría optar por una estrategia conservadora que minimice el riesgo de efectos adversos.

- Criterio de Laplace: Este criterio asume que todos los estados de la naturaleza son igualmente probables y, por lo tanto, se basa en la maximización de la utilidad media de las consecuencias posibles. Esto puede aplicarse en nutrición cuando no se tiene información precisa sobre las probabilidades de diferentes resultados, como cuando se introduce una nueva dieta experimental.
- Criterio de minimax: Se selecciona la acción que minimiza la pérdida máxima posible. En estudios nutricionales, este enfoque puede aplicarse cuando se desean evitar los peores resultados posibles, como recomendar dietas que minimicen el riesgo de malnutrición o deficiencias críticas en una población.

### 1.6.3 Árboles de decisión

Un árbol de decisión es una herramienta gráfica que ayuda a visualizar los diferentes cursos de acción posibles, los estados de la naturaleza y las consecuencias de cada decisión. Los árboles de decisión son útiles para estructurar problemas complejos y tomar decisiones basadas en la comparación de los posibles resultados de diferentes opciones. En nutrición, un árbol de decisión puede ser utilizado para evaluar la mejor intervención dietética considerando diferentes grupos de población, factores de riesgo y resultados esperados.

Por ejemplo, si se está evaluando la implementación de un programa de suplementación de vitamina D, el árbol de decisión podría incluir diferentes opciones de dosis, la probabilidad de deficiencia en diferentes grupos de edad y las posibles consecuencias de la suplementación o falta de la misma.

### 1.6.4 Aplicaciones de la teoría de decisiones en nutrición

En el ámbito de la nutrición, la teoría de decisiones tiene aplicaciones clave, ya que permite a los nutricionistas y profesionales de la salud pública realizar recomendaciones informadas. Algunos ejemplos de aplicación incluyen:

- Decidir si introducir o no un nuevo programa de educación alimentaria en una comunidad, basado en la evidencia previa y las posibles consecuencias.
- Evaluar la efectividad de diferentes dietas en la prevención de enfermedades crónicas, considerando tanto los beneficios como los riesgos asociados a cada dieta.
- Tomar decisiones sobre la administración de suplementos nutricionales en poblaciones vulnerables, donde los datos disponibles pueden ser limitados y las consecuencias de una decisión incorrecta pueden ser significativas.

La teoría de decisiones en estadística proporciona un marco valioso para tomar decisiones racionales y fundamentadas en situaciones de incertidumbre. En el campo de la nutrición, esta teoría permite a los profesionales evaluar diferentes cursos de acción basándose en datos estadísticos, probabilidades y utilidades asociadas con las consecuencias de las decisiones. Con la ayuda de herramientas como los criterios de decisión y los árboles de decisión, los

nutricionistas pueden hacer recomendaciones más precisas y seguras, asegurando que las intervenciones en salud alimentaria sean lo más efectivas posibles para mejorar el bienestar de la población.

## 1.7 Componentes de una investigación estadística

Una investigación estadística en nutrición consta de los siguientes componentes:

- Planteamiento del problema: Definir una pregunta de investigación, como “¿Influye la dieta mediterránea en la reducción de la inflamación en pacientes con enfermedades crónicas?”
- Recolección de datos: Recopilar datos sobre la dieta de los participantes y sus niveles de marcadores inflamatorios.
- Análisis de datos: Utilizar técnicas estadísticas para evaluar la relación entre la dieta y los marcadores de inflamación.
- Conclusiones: A partir de los resultados obtenidos, hacer inferencias sobre la efectividad de la dieta mediterránea en la reducción de la inflamación.

La investigación estadística es un proceso sistemático que busca recolectar, analizar e interpretar datos para resolver preguntas o problemas específicos. Este proceso implica varios componentes que son esenciales para garantizar que los resultados sean válidos y útiles. A continuación, se detallan los principales componentes de una investigación estadística, con un enfoque en su aplicación en el campo de la nutrición.

### 1.7.1 Planteamiento del problema

El primer paso en cualquier investigación estadística es definir claramente el problema de investigación. Este paso implica identificar la pregunta que se desea responder y determinar la relevancia del estudio. En el contexto de la nutrición, esto podría incluir preguntas sobre el impacto de una dieta específica en la salud de una población, la prevalencia de deficiencias nutricionales o la efectividad de un programa educativo sobre hábitos alimenticios.

### 1.7.2 Revisión de la literatura

Antes de avanzar con la investigación, es fundamental realizar una revisión de la literatura existente relacionada con el tema. Este componente implica examinar estudios previos, teorías y enfoques metodológicos relevantes que puedan informar la investigación actual. En nutrición, esto puede incluir la revisión de artículos sobre intervenciones dietéticas, estudios sobre la ingesta de nutrientes, y guías de prácticas recomendadas.

### 1.7.3 Diseño de la investigación

El diseño de la investigación es la estructura que guiará el proceso de recolección y análisis de datos. Existen diferentes tipos de diseños, como estudios descriptivos, experimentales, y observacionales. La elección del diseño depende de la naturaleza del problema y de las preguntas de investigación. Por ejemplo, un estudio experimental podría investigar el efecto de un nuevo suplemento nutricional en una población específica, mientras que un estudio observacional podría analizar patrones de consumo de alimentos en diferentes grupos demográficos.

#### 1.7.4 Población y muestra

En cualquier investigación estadística, es fundamental definir la población de interés, que es el conjunto completo de individuos o elementos sobre los que se desea obtener conclusiones. Dado que a menudo es impracticable estudiar a toda la población, se selecciona una muestra representativa. La forma en que se selecciona la muestra puede afectar la validez de los resultados, por lo que es esencial utilizar métodos de muestreo adecuados. En nutrición, por ejemplo, la muestra podría consistir en un grupo de individuos de diferentes edades y antecedentes culturales para evaluar su ingesta dietética.

#### 1.7.5 Recolección de datos

La recolección de datos implica obtener información sobre la muestra seleccionada. Este proceso puede realizarse a través de encuestas, entrevistas, análisis de registros, o mediciones directas. En el ámbito de la nutrición, los métodos de recolección de datos pueden incluir cuestionarios sobre hábitos alimenticios, registros de consumo de alimentos, y mediciones de indicadores de salud como el índice de masa corporal (IMC) o niveles de nutrientes en sangre.

#### 1.7.6 Análisis de datos

Una vez que se han recolectado los datos, el siguiente paso es el análisis de datos. Este proceso implica utilizar técnicas estadísticas para examinar la información y extraer conclusiones. Dependiendo de la naturaleza de los datos y de las preguntas de investigación, se pueden utilizar análisis descriptivos, inferenciales, o multivariados. En nutrición, el análisis de datos podría ayudar a identificar tendencias en la ingesta de alimentos, correlaciones entre hábitos alimenticios y problemas de salud, o la eficacia de intervenciones nutricionales.

#### 1.7.7 Interpretación de resultados

La interpretación de resultados es un componente crítico de la investigación estadística. Este paso implica evaluar lo que significan los hallazgos y cómo se relacionan con el problema de investigación original. En nutrición, la interpretación puede incluir la identificación de implicaciones para la salud pública, recomendaciones para prácticas dietéticas, o la necesidad de cambios en políticas alimentarias.

#### 1.7.8 Comunicación de resultados

Finalmente, es esencial comunicar los resultados de manera clara y efectiva. Esto puede hacerse a través de informes escritos, presentaciones, o publicaciones en revistas científicas. En el campo de la nutrición, la comunicación de resultados es vital para informar a otros profesionales de la salud, responsables de políticas, y el público en general sobre los hallazgos relevantes que pueden afectar la salud y el bienestar de la población.

Los componentes de una investigación estadística son interdependientes y, al ser implementados de manera efectiva, contribuyen a la calidad y utilidad de la investigación. En el ámbito de la nutrición, estos componentes son fundamentales para abordar preguntas

complejas sobre la relación entre la dieta, la salud y el bienestar, y para guiar la toma de decisiones informadas en la práctica y la política alimentaria.

## 1.8 Recolección de datos

La recolección de datos es uno de los aspectos más importantes en una investigación nutricional. Puede realizarse mediante encuestas dietéticas, diarios alimentarios, o mediciones biométricas. La calidad de los datos recolectados es esencial para asegurar la validez del estudio. En estudios de nutrición, es común utilizar herramientas como el cuestionario de frecuencia alimentaria (CFA) para evaluar la ingesta de alimentos en grandes poblaciones.

La recolección de datos es un componente fundamental en cualquier investigación estadística, ya que la calidad de los datos recopilados influye directamente en la validez y fiabilidad de los resultados. Este proceso implica la obtención de información relevante que permita responder a las preguntas de investigación planteadas. En el contexto de la nutrición, la recolección de datos puede adoptar diversas formas y métodos, cada uno con sus ventajas y desventajas.

### 1.8.1 Fuentes de datos

La recolección de datos puede provenir de diferentes fuentes de datos, que generalmente se dividen en dos categorías:

- **Datos primarios:** Son aquellos datos que se obtienen directamente de la fuente a través de métodos de recolección específicos. Ejemplos de datos primarios en nutrición incluyen encuestas sobre hábitos alimenticios, registros de consumo de alimentos, y mediciones antropométricas (como peso y altura).
- **Datos secundarios:** Son datos que han sido recolectados previamente por otros investigadores o instituciones. Estos pueden incluir estadísticas de salud pública, estudios anteriores sobre nutrición, y registros administrativos. Si bien los datos secundarios pueden ser más fáciles y rápidos de obtener, su calidad y relevancia para el estudio actual deben ser evaluadas cuidadosamente.

### 1.8.2 Métodos de recolección de datos

Existen varios métodos de recolección de datos que se pueden utilizar en investigaciones nutricionales. Algunos de los más comunes incluyen:

- **Encuestas y cuestionarios:** Estos instrumentos permiten recopilar información de manera estructurada y sistemática. Pueden ser administrados en persona, por teléfono, o en línea. Las preguntas pueden ser cerradas (opciones predefinidas) o abiertas (respuestas libres). Las encuestas sobre hábitos alimenticios suelen incluir preguntas sobre frecuencia de consumo de alimentos, preferencias dietéticas, y conocimiento sobre nutrición.
- **Diarios de alimentos:** Este método implica que los participantes registren lo que comen y beben durante un período determinado. Los diarios de alimentos proporcionan datos detallados sobre la ingesta dietética, pero requieren un alto grado de compromiso y precisión por parte de los participantes.

- Entrevistas: Las entrevistas permiten obtener información más profunda y cualitativa sobre los hábitos alimenticios y actitudes hacia la nutrición. Pueden ser estructuradas, semi-estructuradas o no estructuradas, lo que permite flexibilidad en la conversación.
- Observación directa: Este método implica observar directamente los comportamientos de alimentación de los participantes en un entorno natural. La observación puede proporcionar información valiosa sobre la selección de alimentos y los patrones de consumo, aunque puede ser subjetiva y menos cuantificable.
- Mediciones físicas: Incluir la recolección de datos antropométricos (como peso, altura, IMC) y análisis de laboratorio (como niveles de nutrientes en sangre) es esencial para evaluar el estado nutricional de los individuos.

### 1.8.3 Consideraciones éticas

La recolección de datos en investigaciones nutricionales debe llevarse a cabo de manera ética. Esto implica:

- Consentimiento informado: Los participantes deben ser informados sobre el propósito del estudio, los procedimientos, y los posibles riesgos y beneficios antes de consentir su participación.
- Confidencialidad: Es fundamental proteger la identidad y la información personal de los participantes, asegurando que los datos se manejen de manera confidencial.
- Derecho a la retirada: Los participantes deben tener la libertad de retirarse del estudio en cualquier momento sin ninguna repercusión.

### 1.8.4 Validación de datos

Una vez recolectados, los datos deben ser validados para asegurar su precisión y consistencia. Esto puede incluir:

- Verificación de datos: Revisar los registros para detectar errores o inconsistencias.
- Pruebas de confiabilidad: Evaluar si las mediciones y los instrumentos utilizados producen resultados consistentes y fiables a lo largo del tiempo.
- Análisis piloto: Realizar una prueba preliminar del método de recolección de datos en una muestra pequeña para identificar problemas potenciales antes de la recolección a gran escala.

La recolección de datos es un proceso crucial en la investigación estadística que requiere una planificación cuidadosa y un enfoque ético. En el campo de la nutrición, la calidad de los datos recopilados influye directamente en la capacidad de los investigadores para hacer recomendaciones informadas sobre la salud y el bienestar de la población. Una adecuada recolección de datos proporciona la base necesaria para análisis posteriores y para la toma de decisiones fundamentadas en la práctica de la nutrición.

## 1.9 Estadística paramétrica

La estadística paramétrica se basa en suposiciones sobre la distribución de los datos, generalmente asumiendo que siguen una distribución normal. En nutrición, los métodos paramétricos se utilizan frecuentemente para comparar los efectos de diferentes dietas. Por ejemplo, podríamos usar una prueba t para comparar el IMC de personas que siguen una

dieta baja en carbohidratos con las que siguen una dieta alta en carbohidratos.

La estadística paramétrica es un enfoque de análisis estadístico que se basa en suposiciones sobre la distribución de los datos. Este enfoque es fundamental en la inferencia estadística, ya que permite realizar estimaciones y pruebas de hipótesis sobre poblaciones a partir de muestras. En esta sección, se abordarán los conceptos clave, las condiciones de aplicabilidad, y algunos ejemplos de técnicas paramétricas en el contexto de la nutrición.

### 1.9.1 Fundamentos de la estadística paramétrica

La estadística paramétrica se basa en la idea de que los datos provienen de una distribución específica, típicamente la distribución normal. Algunas de las características fundamentales de la estadística paramétrica incluyen:

- **Suposiciones sobre la distribución:** Para que los métodos paramétricos sean válidos, se requiere que los datos cumplan ciertas condiciones, como la normalidad y la homogeneidad de varianzas. Esto significa que la variable que se estudia debe seguir una distribución normal, y las varianzas de los diferentes grupos comparados deben ser aproximadamente iguales.
- **Uso de parámetros poblacionales:** La estadística paramétrica se centra en estimar parámetros poblacionales (como la media y la desviación estándar) a partir de estadísticas muestrales. Las inferencias se realizan utilizando estos parámetros para generalizar a la población.
- **Mayor poder estadístico:** Cuando las suposiciones se cumplen, los métodos paramétricos suelen tener más poder para detectar diferencias y efectos que los métodos no paramétricos, lo que significa que son más eficaces para identificar relaciones significativas entre variables.

### 1.9.2 Condiciones de aplicabilidad

Antes de aplicar métodos paramétricos, es crucial evaluar si los datos cumplen con las suposiciones requeridas. Las principales condiciones incluyen:

- **Normalidad:** Los datos deben aproximarse a una distribución normal. Esta suposición se puede evaluar visualmente a través de gráficos de histograma o QQ plots, y mediante pruebas de normalidad como la prueba de Shapiro-Wilk.
- **Homogeneidad de varianzas:** Al comparar múltiples grupos, es esencial que las varianzas sean aproximadamente iguales. Esto se puede verificar mediante la prueba de Levene o la prueba de Bartlett.
- **Independencia:** Las observaciones deben ser independientes entre sí, lo que significa que la medición de un individuo no debe influir en la medición de otro.

### 1.9.3 Ejemplos de métodos estadísticos paramétricos

En el ámbito de la nutrición, varios métodos estadísticos paramétricos son comúnmente utilizados. Algunos ejemplos incluyen:

- **Prueba t de Student:** Utilizada para comparar las medias de dos grupos independientes (por ejemplo, comparar la ingesta calórica entre dos poblaciones distintas). La prueba

t determina si las diferencias observadas son estadísticamente significativas.

- ANOVA (Análisis de Varianza): Se utiliza para comparar las medias de tres o más grupos. Por ejemplo, se podría utilizar ANOVA para analizar si hay diferencias significativas en el índice de masa corporal (IMC) entre diferentes grupos de edad que siguen distintas dietas.
- Regresión lineal: Este método se emplea para modelar la relación entre una variable dependiente y una o más variables independientes. En nutrición, la regresión lineal podría utilizarse para analizar cómo la ingesta de ciertos nutrientes afecta el peso corporal.
- Correlación de Pearson: Se utiliza para medir la fuerza y dirección de la relación lineal entre dos variables cuantitativas. Por ejemplo, se puede analizar la relación entre la ingesta de frutas y verduras y los niveles de colesterol en sangre.

#### 1.9.4 Ventajas y desventajas

La estadística paramétrica presenta varias ventajas y desventajas:

Ventajas

- Proporciona estimaciones precisas y fiables cuando se cumplen las suposiciones.
- Permite realizar inferencias más potentes y detalladas sobre la población.
- Los resultados son generalmente más fáciles de interpretar y comunicar.

Desventajas

- Puede dar lugar a resultados engañosos si las suposiciones no se cumplen.
- Menos robusta en situaciones donde los datos son escasos o no siguen una distribución normal.
- Puede ser inapropiada para datos categóricos o de ordinalidad.

La estadística paramétrica es una herramienta valiosa en la investigación nutricional que permite realizar análisis significativos bajo ciertas condiciones. Comprender las suposiciones y aplicaciones de estos métodos es crucial para obtener conclusiones válidas y aplicables en el ámbito de la salud y la nutrición. El uso adecuado de la estadística paramétrica puede contribuir significativamente a la toma de decisiones informadas en políticas y prácticas alimentarias.

#### 1.10 Población

En estadística, la población se refiere al conjunto completo de individuos o elementos sobre los cuales se desea hacer inferencias. En estudios de nutrición, la población puede ser, por ejemplo, todos los adultos de un país, o todos los pacientes con diabetes tipo 2 en una clínica.

En estadística, el término población se refiere al conjunto total de elementos que comparten una característica específica y que son objeto de estudio. Comprender el concepto de población es fundamental para el desarrollo de investigaciones estadísticas, ya que todas las conclusiones e inferencias se derivan de los datos recolectados a partir de esta. En esta sección, exploraremos las definiciones, características y tipos de poblaciones, así como su importancia en el contexto de la investigación nutricional.

### 1.10.1 Definición de población

La población se define como el conjunto completo de individuos, objetos o eventos que comparten una o más características comunes. Por ejemplo, si se está estudiando la ingesta de nutrientes en adultos de una determinada región, la población incluiría a todos los adultos en esa región.

### 1.10.2 Características de la población

Las poblaciones pueden caracterizarse por varias propiedades importantes:

- **Tamaño de la población:** Se refiere a la cantidad total de elementos que conforman la población. Puede ser finita (un número limitado de individuos) o infinita (individuos en constante cambio o expansión).
- **Homogeneidad y heterogeneidad:** Una población puede ser homogénea si todos los elementos son similares en relación con la característica que se está estudiando, o heterogénea si hay variaciones significativas entre los elementos.
- **Estructura:** La estructura de una población se puede analizar en términos de su composición demográfica (edad, género, raza, etc.) y otras características relevantes (nivel socioeconómico, hábitos alimenticios, etc.).

### 1.10.3 Tipos de poblaciones

Las poblaciones pueden clasificarse de diversas maneras:

- **Poblaciones finitas e infinitas:** Como se mencionó anteriormente, una población finita tiene un número determinado de elementos, mientras que una población infinita no tiene un límite definido.
- **Poblaciones homogéneas y heterogéneas:** Las poblaciones homogéneas presentan características similares, lo que facilita el análisis, mientras que las poblaciones heterogéneas pueden requerir técnicas más complejas para obtener conclusiones significativas.
- **Poblaciones de estudio y poblaciones de referencia:** La población de estudio es el grupo específico que se investiga, mientras que la población de referencia puede ser un grupo más amplio que sirve como comparación.

### 1.10.4 Importancia de la población en la investigación nutricional

Entender la población es crucial para la investigación en nutrición por varias razones:

- **Definición de objetivos:** La identificación clara de la población objetivo ayuda a establecer los objetivos de la investigación, lo que permite formular preguntas de investigación relevantes y específicas.
- **Selección de la muestra:** Conocer las características de la población es esencial para seleccionar una muestra representativa, lo que a su vez garantiza que los resultados puedan generalizarse a la población más amplia.
- **Interpretación de resultados:** Los hallazgos de una investigación solo son útiles si se comprenden dentro del contexto de la población estudiada. La interpretación adecuada de los resultados depende de cómo se relacionan con las características de la población.

### 1.10.5 Ejemplos de poblaciones en nutrición

En el ámbito de la nutrición, las poblaciones pueden variar ampliamente:

- Una población puede ser todos los adolescentes de una ciudad que están en riesgo de obesidad debido a malos hábitos alimenticios.
- Otra población podría incluir a todas las mujeres embarazadas en un país que consumen suplementos vitamínicos.
- También se podría estudiar la población de ancianos que participan en un programa de nutrición diseñado para mejorar su salud.

El concepto de población es un pilar fundamental en la investigación estadística. Al comprender las características, tipos y la importancia de la población en el contexto de la nutrición, los investigadores pueden llevar a cabo estudios más precisos y significativos. La identificación clara de la población y sus características permite establecer objetivos claros, seleccionar muestras representativas y obtener resultados que tengan un impacto real en la salud y el bienestar de la comunidad.

## 1.11 Muestra aleatoria

Una muestra aleatoria es un subconjunto de la población que se selecciona de manera que cada individuo tiene la misma probabilidad de ser elegido. En nutrición, trabajar con muestras aleatorias es crucial para evitar sesgos y asegurar que los resultados obtenidos en el estudio sean representativos de la población general. Por ejemplo, para estudiar la ingesta de proteínas en la población mexicana, se selecciona una muestra aleatoria de individuos de diversas regiones, edades y estilos de vida.

El concepto de muestra aleatoria es fundamental en la estadística, ya que se refiere a un subconjunto de elementos extraídos de una población de manera que cada elemento tenga la misma probabilidad de ser seleccionado. Esta técnica garantiza que la muestra sea representativa de la población, lo que permite realizar inferencias y generalizaciones válidas. En esta sección, exploraremos la definición de muestra aleatoria, sus tipos, métodos de selección, importancia y ejemplos en el contexto de la nutrición.

### 1.11.1 Definición de muestra aleatoria

Una muestra aleatoria se define como un subconjunto de individuos seleccionado de una población de manera que cada miembro de la población tenga una probabilidad conocida y no cero de ser elegido. Este enfoque minimiza el sesgo y permite que los resultados obtenidos sean generalizables a la población en su totalidad.

### 1.11.2 Importancia de la muestra aleatoria

La selección de una muestra aleatoria es crucial por varias razones:

- Representatividad: Asegura que la muestra refleje las características de la población, lo que es esencial para la validez de los resultados.
- Reducción del sesgo: Al eliminar la influencia de factores subjetivos en la selección de la muestra, se reduce el riesgo de sesgo sistemático.

- Generalización: Permite hacer inferencias sobre la población en base a los resultados obtenidos de la muestra, lo que es fundamental en estudios de nutrición y salud pública.

### 1.11.3 Tipos de muestreo aleatorio

Existen varios tipos de muestreo aleatorio, cada uno con sus propias características y aplicaciones:

- Muestreo aleatorio simple: Cada individuo en la población tiene la misma probabilidad de ser seleccionado. Este tipo de muestreo se puede realizar utilizando un generador de números aleatorios o sorteos.
- Muestreo sistemático: Se selecciona un punto de partida aleatorio y se elige cada  $k$ -ésimo elemento de la población. Por ejemplo, si se tiene una lista de 1000 individuos y se quiere una muestra de 100, se podría seleccionar cada décimo individuo.
- Muestreo estratificado: La población se divide en grupos (estratos) basados en características relevantes (por ejemplo, edad, género) y se selecciona una muestra aleatoria de cada estrato. Este método es útil cuando se quiere asegurar la representación de subgrupos específicos dentro de la población.
- Muestreo por conglomerados: Se divide la población en grupos (conglomerados) y se selecciona aleatoriamente algunos de estos grupos. Luego, se estudian todos los elementos dentro de los grupos seleccionados. Este método es eficiente y práctico, especialmente en investigaciones a gran escala.

### 1.11.4 Métodos de selección de muestra aleatoria

La elección del método de muestreo depende de la naturaleza de la población y los objetivos de la investigación. Algunos métodos incluyen:

- Listas de población: Utilizar listas actualizadas de la población para asegurar que todos los individuos tienen la misma oportunidad de ser seleccionados.
- Software de muestreo: Utilizar programas informáticos que generan muestras aleatorias para garantizar la imparcialidad.
- Diseños de encuesta: Implementar diseños de encuesta que aseguren la aleatoriedad en la selección de participantes.

### 1.11.5 Ejemplos de muestreo aleatorio en nutrición

En el contexto de la nutrición, el muestreo aleatorio se aplica en diversas investigaciones:

- Un estudio que busca evaluar la ingesta de frutas y verduras en una población de adultos puede utilizar muestreo aleatorio simple para seleccionar participantes de una lista de residentes de una ciudad.
- En un programa de intervención nutricional dirigido a adolescentes, se puede aplicar muestreo estratificado para asegurar que se incluya una representación adecuada de diferentes grupos etarios y socioeconómicos.
- En una investigación sobre hábitos alimentarios en diferentes regiones de un país, el muestreo por conglomerados puede ser útil para seleccionar aleatoriamente ciertas ciudades y estudiar a todos los residentes de esas áreas.

---

La selección de una muestra aleatoria es un componente esencial de la investigación estadística, ya que garantiza la representatividad y la validez de los resultados. En el ámbito de la nutrición, el uso de muestreo aleatorio permite a los investigadores hacer generalizaciones confiables sobre hábitos alimentarios y su impacto en la salud. Al comprender los diferentes tipos de muestreo y sus aplicaciones, los investigadores pueden diseñar estudios que ofrezcan hallazgos significativos y aplicables.



# INFERENCIA ESTADÍSTICA: ESTIMACIÓN, MUESTREO

## 2 — Inferencia Estadística: Estimación, Muestreo

### 2.1 Teoría de conjuntos

La teoría de conjuntos es un área fundamental en matemáticas que se utiliza para describir y analizar grupos de objetos o elementos. En el contexto de la estadística, se utiliza para organizar y categorizar datos. A continuación se presentan algunos conceptos clave relacionados con la teoría de conjuntos:

#### 2.1.1 Conjuntos

Un conjunto es una colección de elementos que comparten una característica común. Se denota comúnmente con letras mayúsculas y se representa entre llaves. Por ejemplo, el conjunto  $A = \{1, 2, 3, 4, 5\}$  contiene los números del uno al cinco.

#### 2.1.2 Operaciones con conjuntos

Las operaciones más comunes en teoría de conjuntos incluyen:

- Unión ( $A \cup B$ ): El conjunto de elementos que están en  $A$ , en  $B$  o en ambos.
- Intersección ( $A \cap B$ ): El conjunto de elementos que están en ambos conjuntos  $A$  y  $B$ .
- Diferencia ( $A - B$ ): El conjunto de elementos que están en  $A$  pero no en  $B$ .
- Complemento ( $A'$ ): El conjunto de elementos que no están en  $A$ .

#### 2.1.3 Diagramas de Venn

Los diagramas de Venn son representaciones gráficas que muestran las relaciones entre conjuntos. Estos diagramas son útiles para visualizar la unión, intersección y diferencia de conjuntos. Por ejemplo, el siguiente diagrama muestra tres conjuntos  $P$ ,  $C$  y  $G$ :



#### 2.1.4 Definición de Conjunto

Un conjunto se define como una colección de elementos distintos, considerados como un objeto en sí mismo. Los conjuntos se denotan comúnmente con letras mayúsculas, y sus elementos se encierran entre llaves. Por ejemplo, el conjunto de números naturales menores que 5 se puede representar como:

$$A = \{0, 1, 2, 3, 4\}$$

#### 2.1.5 Tipos de Conjuntos

Los conjuntos pueden clasificarse de varias maneras:

- Conjunto vacío: Es el conjunto que no contiene elementos, denotado por  $\emptyset$  o  $\{\}$ .
- Conjunto finito: Un conjunto que contiene un número limitado de elementos, como  $B = \{1, 2, 3\}$ .
- Conjunto infinito: Un conjunto que tiene un número ilimitado de elementos, como el conjunto de todos los números naturales  $\mathbb{N}$ .
- Conjuntos universales: Es el conjunto que contiene todos los posibles elementos bajo consideración en un contexto específico, denotado comúnmente por  $U$ .

#### 2.1.6 Notación de Conjuntos

Existen diversas formas de representar los elementos de un conjunto, tales como:

- Listando los elementos: Como en el ejemplo anterior  $A = \{0, 1, 2, 3, 4\}$ .
- Definición por propiedades: En lugar de listar todos los elementos, se define el conjunto por una propiedad común. Por ejemplo, el conjunto de números pares puede escribirse como:

$$C = \{x \in \mathbb{Z} \mid x \text{ es par}\}$$

### 2.1.7 Operaciones con Conjuntos

Las operaciones con conjuntos permiten combinar y relacionar conjuntos de diferentes maneras. Las principales operaciones son:

- Unión ( $A \cup B$ ): Es el conjunto de elementos que están en  $A$ , en  $B$  o en ambos.
- Intersección ( $A \cap B$ ): Es el conjunto de elementos que están en ambos conjuntos  $A$  y  $B$ .
- Diferencia ( $A - B$ ): Es el conjunto de elementos que están en  $A$  pero no en  $B$ .
- Complemento ( $A'$ ): Es el conjunto de todos los elementos que no están en  $A$  respecto al conjunto universal  $U$ .

### 2.1.8 Propiedades de los Conjuntos

Existen diversas propiedades que los conjuntos pueden satisfacer, entre ellas:

- Conmutatividad:

$$A \cup B = B \cup A$$

$$A \cap B = B \cap A$$

- Asociatividad:

$$(A \cup B) \cup C = A \cup (B \cup C)$$

$$(A \cap B) \cap C = A \cap (B \cap C)$$

- Distributividad:

$$A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$$

$$A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$$

### 2.1.9 Aplicaciones en Estadística

La teoría de conjuntos es fundamental para la inferencia estadística, ya que permite clasificar y analizar datos de manera estructurada, definir claramente poblaciones, muestras y eventos dentro de experimentos estadísticos. Los conceptos de conjuntos ayudan a organizar los datos y establecer relaciones entre diferentes grupos, facilitando así el análisis estadístico.

La teoría de conjuntos es una rama fundamental de las matemáticas que se ocupa del estudio de colecciones de objetos, llamados elementos o miembros. En estadística, esta teoría es esencial para la organización y análisis de datos. A continuación se presentan los conceptos clave relacionados con la teoría de conjuntos.

## 2.2 Distribución de Muestreo

La distribución de muestreo es un concepto fundamental en la estadística inferencial que describe cómo varían las estadísticas muestrales al tomar múltiples muestras de una misma población. Esta distribución es crucial para entender cómo los estimadores, como la media

o la proporción, se comportan en relación con la población de la que se extraen, especialmente en estudios relacionados con la nutrición y la salud.

La distribución de muestreo se refiere a la distribución de un estadístico (por ejemplo, la media de consumo calórico diario) que se calcula a partir de todas las posibles muestras de un tamaño específico tomadas de una población. Esta distribución proporciona información sobre la variabilidad y el comportamiento del estadístico en diferentes muestras. Por lo general, se puede definir para cualquier estadístico, pero las más comunes son la media muestral y la proporción muestral.

### 2.2.1 Definición

La distribución de muestreo es esencial para estimar parámetros poblacionales, como el consumo promedio de nutrientes (carbohidratos, proteínas, grasas) entre diferentes grupos poblacionales. Al calcular el consumo de estos nutrientes en diferentes muestras de la población, podemos obtener una idea precisa de la ingesta nutricional general y su variabilidad.

### 2.2.2 Teorema del Límite Central

El Teorema Central del Límite establece que, bajo ciertas condiciones, la distribución de la media de una muestra tomará una forma aproximadamente normal a medida que el tamaño de la muestra aumenta, independientemente de la forma de la distribución de la población original. Esto es fundamental para la inferencia estadística, ya que permite utilizar la distribución normal para realizar estimaciones y pruebas de hipótesis en estudios de nutrición.

Este teorema implica que:

- La media de la distribución de muestreo de las medias muestrales es igual a la media de la población ( $\mu$ ), que podría ser, por ejemplo, la ingesta calórica promedio de la población objetivo.
- La desviación estándar de la distribución de muestreo de las medias muestrales, conocida como error estándar ( $\sigma_{\bar{x}}$ ), se calcula como:

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

donde  $\sigma$  es la desviación estándar de la población de ingesta nutricional y  $n$  es el tamaño de la muestra.

### 2.2.3 Importancia de la Distribución de Muestreo

La comprensión de la distribución de muestreo es esencial por varias razones en el ámbito de la nutrición:

- Permite a los investigadores hacer inferencias sobre la ingesta nutricional de la población a partir de muestras.
- Ayuda a determinar la precisión de un estimador, lo que se refleja en la amplitud del intervalo de confianza. Esto es vital para evaluar si las recomendaciones dietéticas son efectivas.

- Facilita la realización de pruebas de hipótesis, permitiendo la comparación de estadísticas muestrales de grupos con diferentes hábitos alimenticios o condiciones de salud.

#### 2.2.4 Ejemplo de Distribución de Muestreo

Supongamos que tenemos una población de  $N = 1000$  adultos y estamos interesados en la media de su consumo diario de calorías. Si tomamos múltiples muestras de tamaño  $n = 30$  y calculamos la media de cada muestra, estas medias muestrales formarán una nueva distribución. A medida que aumentamos el número de muestras, la distribución de estas medias muestrales se acercará a una distribución normal, según lo predicho por el Teorema Central del Límite.

#### 2.2.5 Distribuciones de Muestreo Comunes

Existen varias distribuciones de muestreo que se utilizan con frecuencia en estudios de nutrición:

- Distribución de la Media Muestral: Distribución de las medias calculadas a partir de múltiples muestras de ingesta de nutrientes.
- Distribución de la Proporción Muestral: Distribución de las proporciones de la población que cumplen con ciertos estándares nutricionales (por ejemplo, el porcentaje de personas que cumplen con las recomendaciones de ingesta de frutas y verduras).
- Distribución de la Varianza Muestral: Distribución de las varianzas de la ingesta de nutrientes calculadas a partir de múltiples muestras, relevante para análisis de varianza en estudios de dieta.

#### 2.2.6 Consideraciones sobre la Muestra

Es importante considerar el tamaño de la muestra y la técnica de muestreo utilizada en estudios de nutrición. El tamaño de la muestra  $n$  debe ser suficientemente grande para que la distribución de muestreo se asemeje a la normal. Para poblaciones no normales, generalmente se recomienda un tamaño de muestra de al menos 30. Además, se debe aplicar un muestreo aleatorio para garantizar que la muestra sea representativa de la población, lo que es crucial para obtener conclusiones válidas sobre la ingesta nutricional.

La distribución de muestreo es un concepto clave en la estadística inferencial, ya que permite realizar inferencias sobre una población basándose en muestras. A través del Teorema Central del Límite, se puede comprender cómo las medias muestrales tienden a seguir una distribución normal a medida que se aumenta el tamaño de la muestra, lo que es fundamental para el análisis y la interpretación de datos en estudios estadísticos, especialmente en el ámbito de la nutrición.

### 2.3 Muestreo Aleatorio Simple

El muestreo aleatorio simple es una técnica fundamental en la estadística que permite seleccionar un subconjunto de individuos de una población de manera que cada individuo tenga la misma probabilidad de ser elegido. Este método es especialmente relevante en estudios

de nutrición, donde se busca obtener información representativa sobre hábitos alimentarios, ingesta calórica, y otros indicadores de salud.

### 2.3.1 Definición

El muestreo aleatorio simple implica la selección de elementos de la población de forma completamente al azar. Cada individuo tiene una probabilidad igual de ser incluido en la muestra, lo que garantiza que la muestra sea representativa de la población en su conjunto. Esto es crucial para evitar sesgos y asegurar que los resultados sean generalizables a la población.

### 2.3.2 Proceso de Muestreo

El proceso de muestreo aleatorio simple generalmente implica los siguientes pasos:

- **Definición de la Población:** Identificar a la población objetivo, que podría ser un grupo de personas, como adultos de una región específica que se estudian para conocer su ingesta nutricional.
- **Determinación del Tamaño de la Muestra:** Decidir cuántos individuos se incluirán en la muestra. Un tamaño de muestra más grande proporciona estimaciones más precisas.
- **Selección Aleatoria:** Utilizar métodos aleatorios, como un generador de números aleatorios, para seleccionar individuos de la lista de la población. Esto puede hacerse mediante la asignación de un número a cada individuo y seleccionando números al azar.

### 2.3.3 Ejemplo en Nutrición

Supongamos que un investigador desea estudiar el consumo de frutas y verduras en una población de 1,000 adultos en una ciudad. Para obtener una muestra representativa, el investigador:

- Define la población como los 1,000 adultos.
- Decide que la muestra tendrá 100 individuos.
- Asigna un número del 1 al 1,000 a cada adulto y utiliza un generador de números aleatorios para seleccionar 100 números. Los adultos correspondientes a esos números serán incluidos en la muestra.

### 2.3.4 Ventajas del Muestreo Aleatorio Simple

El muestreo aleatorio simple ofrece varias ventajas en estudios de nutrición:

- **Simplicidad:** Es un método fácil de entender y de implementar.
- **Representatividad:** Ayuda a garantizar que la muestra sea representativa de la población, lo que mejora la validez de los resultados.
- **Facilidad en el Análisis:** Permite aplicar métodos estadísticos que asumen que los datos son independientes y que la muestra es aleatoria.

### 2.3.5 Desventajas del Muestreo Aleatorio Simple

A pesar de sus ventajas, el muestreo aleatorio simple también tiene algunas desventajas:

- No es Práctico en Poblaciones Muy Grandes: Puede ser difícil y costoso acceder a una lista completa de la población.
- Variabilidad: La variabilidad en las muestras puede ser alta si la población es heterogénea.

El muestreo aleatorio simple es una técnica efectiva para realizar estudios en nutrición, proporcionando una base sólida para hacer inferencias sobre la población. Al garantizar que cada individuo tenga la misma probabilidad de ser seleccionado, este método contribuye a la obtención de datos precisos y representativos, fundamentales para el análisis y la comprensión de los patrones de consumo nutricional.

## 2.4 Muestreo Aleatorio Estratificado Simple

### 2.4.1 Definición

El muestreo aleatorio estratificado simple es una técnica que mejora la precisión de las estimaciones al dividir la población en estratos o subgrupos homogéneos y luego seleccionar una muestra aleatoria simple de cada uno. La división en estratos se realiza en función de una característica de interés que puede influir en la variable que se está estudiando (por ejemplo, edad, sexo, nivel socioeconómico, etc.).

Este método asegura que cada subgrupo esté representado en la muestra, lo que mejora la exactitud de las conclusiones que se puedan hacer sobre la población total.

### 2.4.2 Características del Muestreo Aleatorio Estratificado Simple

- Homogeneidad dentro de los estratos: Cada estrato está compuesto por individuos que son más similares entre sí en relación con la característica que define el estrato.
- Heterogeneidad entre estratos: Los diferentes estratos tienen características que los distinguen claramente entre sí.
- Representatividad: Asegura que todos los subgrupos importantes de la población estén representados en la muestra final.

### 2.4.3 Pasos para implementar el Muestreo Aleatorio Estratificado Simple

1. Identificar la población total.
  - Definir quiénes forman parte de la población que se va a estudiar. Por ejemplo, si queremos estudiar el consumo de calorías en una población, la población podría ser “todos los adultos mayores de 18 años en una ciudad”.
2. Dividir la población en estratos.
  - Seleccionar una o varias características relevantes para dividir la población en subgrupos homogéneos. Un ejemplo podría ser estratificar por género (hombres y mujeres) o por nivel de actividad física (baja, media, alta).
3. Seleccionar una muestra aleatoria de cada estrato.
  - Dentro de cada estrato, realizar un muestreo aleatorio simple para seleccionar a los participantes. Esto asegura que cada grupo esté representado de manera justa en la muestra.

4. Unir las muestras de cada estrato.
  - Combinar las muestras de cada estrato para formar la muestra final.

#### 2.4.4 Ventajas del Muestreo Aleatorio Estratificado Simple

- Mayor precisión en la estimación de parámetros poblacionales: Al reducir la variabilidad dentro de cada estrato, se mejora la precisión de las estimaciones respecto a la población total.
- Asegura la representación adecuada de todos los subgrupos: Esto es crucial en estudios donde algunos grupos pueden estar subrepresentados si solo se utiliza un muestreo aleatorio simple.
- Permite realizar análisis detallados por estrato: Al separar la muestra por subgrupos, es posible analizar cómo se comportan los diferentes estratos.

#### 2.4.5 Desventajas

- Mayor complejidad logística: Requiere conocimiento previo de la población para poder identificar y dividirla en estratos. Esto puede ser más costoso y difícil de implementar que el muestreo aleatorio simple.
- No es útil si los estratos no son significativamente diferentes: Si los subgrupos no presentan diferencias importantes, dividir la población en estratos puede no aportar beneficios.

#### 2.4.6 Consideraciones

- Proporcionalidad: Es importante decidir si se tomará una muestra proporcional de cada estrato (estratificado proporcional) o si se usará un tamaño de muestra fijo para cada estrato (estratificado no proporcional). En el caso del estratificado proporcional, la muestra de cada estrato será proporcional a su tamaño en la población.
- Fórmula para estimación del tamaño de muestra en un muestreo estratificado simple: Para determinar el tamaño de la muestra en cada estrato, podemos utilizar la fórmula:

$$n_h = \frac{N_h}{N} \cdot n$$

Donde:

- $n_h$  es el tamaño de la muestra en el estrato  $h$ ,
- $N_h$  es el tamaño del estrato  $h$  en la población total,
- $N$  es el tamaño de la población total,
- $n$  es el tamaño total de la muestra que se desea obtener.

#### 2.4.7 Ejemplo de Muestreo Aleatorio Estratificado Simple

Supongamos que queremos realizar un estudio sobre la ingesta calórica de los empleados en una empresa que tiene un total de  $N = 1,000$  empleados.

La empresa tiene tres departamentos:

- Producción:  $N_h = 600$  empleados
- Ventas:  $N_h = 300$  empleados

- Administración:  $N_h = 100$  empleados

Decidimos que queremos una muestra total de  $n = 200$  empleados para el estudio.

Objetivo

Queremos determinar cuántos empleados seleccionar de cada departamento (estrato) usando la fórmula para muestreo estratificado.

1. Para el departamento de Producción:

$$n_{\text{producción}} = \frac{N_{\text{producción}}}{N} \cdot n$$

$$n_{\text{producción}} = \frac{600}{1000} \cdot 200 = 0,6 \cdot 200 = 120$$

2. Para el departamento de Ventas:

$$n_{\text{ventas}} = \frac{N_{\text{ventas}}}{N} \cdot n$$

$$n_{\text{ventas}} = \frac{300}{1000} \cdot 200 = 0,3 \cdot 200 = 60$$

3. Para el departamento de Administración:

$$n_{\text{administración}} = \frac{N_{\text{administración}}}{N} \cdot n$$

$$n_{\text{administración}} = \frac{100}{1000} \cdot 200 = 0,1 \cdot 200 = 20$$

Resultados: Por lo tanto, para nuestra muestra de 200 empleados, seleccionaremos:

- 120 empleados del departamento de Producción.
- 60 empleados del departamento de Ventas.
- 20 empleados del departamento de Administración.

Interpretación

Este enfoque garantiza que cada departamento esté representado de manera proporcional en la muestra total. Esto es importante para obtener resultados que reflejen con precisión la ingesta calórica promedio de todos los empleados de la empresa.

#### 2.4.8 Aplicaciones típicas

- Nutrición: Estratificar una población por edad y sexo para evaluar el consumo promedio de calorías o nutrientes en cada subgrupo.
- Encuestas de salud: Estratificar por zonas geográficas o nivel socioeconómico para asegurar que todos los grupos estén representados.
- Educación: Estratificar por nivel educativo o grado académico para asegurar que todos los niveles sean incluidos en el estudio.

El muestreo aleatorio estratificado simple es una técnica eficaz para garantizar que los diferentes subgrupos de una población estén representados en un estudio. Esto es especialmente importante cuando se sabe que existen diferencias entre los subgrupos en relación con la variable de interés. Aunque requiere más información previa sobre la población, su uso aumenta la precisión de los estudios y permite análisis detallados dentro de cada subgrupo.

## 2.5 Muestreo por Conglomerado

### 2.5.1 Definición

El muestreo por conglomerado es una técnica de muestreo en la que la población se divide en grupos o conglomerados naturales, y luego se seleccionan algunos de esos conglomerados al azar para el estudio. Dentro de los conglomerados seleccionados, se estudian todos los individuos o una muestra de ellos. Este método es especialmente útil cuando la población es geográficamente dispersa o se organiza en grupos naturales.

A diferencia del muestreo estratificado, en el muestreo por conglomerado se seleccionan grupos completos de la población, no individuos dentro de cada subgrupo.

### 2.5.2 Características del Muestreo por Conglomerado

- Conglomerados naturales: Los conglomerados son grupos que ya existen de forma natural en la población, como hogares, escuelas, hospitales, vecindarios o comunidades.
- Reducción de costos y tiempo: El muestreo por conglomerado es ideal cuando es costoso o difícil acceder a todos los individuos de la población dispersa.
- Variabilidad entre conglomerados: Este tipo de muestreo puede ser menos eficiente que otros métodos si los conglomerados no son homogéneos internamente o si difieren mucho entre ellos.

### 2.5.3 Tipos de Muestreo por Conglomerado

- Muestreo por conglomerado de una etapa: Se seleccionan al azar varios conglomerados y se estudian todos los individuos dentro de esos conglomerados.
- Muestreo por conglomerado de dos etapas: Primero se seleccionan al azar varios conglomerados, y luego se toma una muestra aleatoria dentro de cada conglomerado seleccionado.
- Muestreo por conglomerado de múltiples etapas: Es una extensión del muestreo de dos etapas, en la que se seleccionan conglomerados en diferentes niveles (por ejemplo, regiones, luego distritos dentro de las regiones, y finalmente escuelas dentro de los distritos).

### 2.5.4 Ventajas del Muestreo por Conglomerado

- Eficiencia en términos de costos y tiempo: Cuando la población está dispersa en un área geográfica extensa, este método reduce los costos de desplazamiento y tiempo al seleccionar grupos de personas en lugar de individuos dispersos.
- Facilidad logística: Es más fácil y práctico estudiar conglomerados completos, como

escuelas o barrios, que estudiar a personas individuales distribuidas por toda la población.

- Aplicable cuando no se tiene un marco de muestreo completo: Si no se tiene una lista completa de todos los individuos de la población, pero sí de los conglomerados, este método facilita el muestreo.

### 2.5.5 Desventajas del Muestreo por Conglomerado

- Menor precisión comparado con el muestreo estratificado: Los conglomerados pueden ser más heterogéneos entre sí, lo que puede llevar a una mayor variabilidad en los resultados. Esto hace que el error estándar del estimador sea mayor que en el muestreo aleatorio simple o estratificado.
- Sesgo si los conglomerados son diferentes: Si hay grandes diferencias entre los conglomerados y no se seleccionan suficientes conglomerados, los resultados pueden no representar bien a la población.
- Dependencia de los elementos dentro de los conglomerados: Los elementos dentro de un conglomerado pueden ser más similares entre sí, lo que disminuye la independencia de las observaciones y puede afectar la validez de los resultados.

### 2.5.6 Pasos para Implementar el Muestreo por Conglomerado

1. Definir los conglomerados: Dividir la población en conglomerados naturales. Por ejemplo, si se desea estudiar los hábitos alimentarios de una ciudad, los conglomerados podrían ser los vecindarios.
2. Seleccionar los conglomerados: Elegir aleatoriamente un número de conglomerados para el estudio. Puede seleccionarse una muestra de conglomerados proporcional al tamaño de cada uno, o usar un enfoque de muestreo de dos etapas para seleccionar individuos dentro de los conglomerados.
3. Recolectar datos de los conglomerados seleccionados: Una vez que se seleccionan los conglomerados, se estudian todos los individuos dentro de los conglomerados (muestreo de una etapa) o se selecciona una muestra dentro de cada conglomerado (muestreo de dos etapas).

### 2.5.7 Ejemplo de Muestreo por Conglomerado

Problema: Supongamos que queremos estudiar la ingesta calórica diaria de los estudiantes de escuelas secundarias en una ciudad. La ciudad tiene 50 escuelas y no es práctico hacer un muestreo individual en todas las escuelas.

Solución usando muestreo por conglomerado:

- Dividir la población en conglomerados: Cada escuela se considera un conglomerado.
- Seleccionar al azar 10 escuelas (conglomerados): De las 50 escuelas, seleccionamos al azar 10 para nuestro estudio.
- Recolectar datos dentro de los conglomerados seleccionados: Podríamos estudiar a todos los estudiantes de las 10 escuelas seleccionadas (muestreo de una etapa) o tomar una muestra aleatoria de estudiantes dentro de cada escuela (muestreo de dos etapas).

Este enfoque reduce significativamente el costo y tiempo del estudio, ya que en lugar de

desplazarnos por todas las escuelas de la ciudad, nos enfocamos en solo 10 escuelas.

### 2.5.8 Aplicaciones Comunes del Muestreo por Conglomerado

- Investigaciones en educación: Las escuelas, aulas o distritos escolares suelen ser conglomerados naturales. Un estudio podría seleccionar algunas escuelas o distritos y luego estudiar a los estudiantes dentro de esos conglomerados.
- Encuestas de hogares: En investigaciones de salud pública o de bienestar económico, las comunidades o vecindarios pueden ser seleccionados como conglomerados. Dentro de los vecindarios seleccionados, se estudian todos los hogares o una muestra de ellos.
- Estudios geográficos: Cuando una población está dispersa geográficamente, se puede dividir en conglomerados según áreas geográficas (por ejemplo, regiones, ciudades o vecindarios), y luego estudiar algunas de esas áreas seleccionadas.

### 2.5.9 Fórmula del Error Estándar en el Muestreo por Conglomerado

El error estándar en el muestreo por conglomerado es mayor que en otros tipos de muestreo debido a la homogeneidad interna de los conglomerados. La fórmula del error estándar puede expresarse como:

$$SE = \sqrt{\frac{S^2}{n} + \frac{S_c^2}{n_c}}$$

Donde:

- $S^2$  es la varianza dentro de los conglomerados,
- $n$  es el número de individuos seleccionados en los conglomerados,
- $S_c^2$  es la varianza entre los conglomerados,
- $n_c$  es el número de conglomerados seleccionados.

Un mayor número de conglomerados seleccionados  $n_c$  disminuye el error estándar, lo que mejora la precisión de los resultados.

### 2.5.10 Consideraciones Finales

El muestreo por conglomerado es una técnica eficiente para estudiar poblaciones grandes y dispersas, ya que reduce costos y facilita la logística de la recolección de datos. Sin embargo, se debe tener cuidado con la posible variabilidad entre conglomerados, que puede afectar la precisión de las estimaciones. Por ello, es recomendable seleccionar un número suficiente de conglomerados para garantizar la representatividad y reducir el error estándar.

## 2.6 Intervalo de Confianza para la Diferencia entre Medias

### 2.6.1 Definición

El intervalo de confianza para la diferencia entre medias es una técnica utilizada en la inferencia estadística para estimar la diferencia entre las medias de dos poblaciones, basándose en muestras de esas poblaciones. Este intervalo proporciona un rango de valores en el que se espera que se encuentre la diferencia entre las medias de las dos poblaciones con un cierto

nivel de confianza.

El intervalo de confianza nos permite decir, por ejemplo, que con un 95% de confianza, la diferencia entre las medias de las dos poblaciones caerá dentro de un rango específico. Si el intervalo incluye el valor cero, esto sugiere que las medias de ambas poblaciones podrían no ser significativamente diferentes.

### 2.6.2 Fórmula

Si tenemos dos muestras independientes con tamaños  $n_1$  y  $n_2$ , medias  $\bar{X}_1$  y  $\bar{X}_2$ , y desviaciones estándar  $s_1$  y  $s_2$ , el intervalo de confianza para la diferencia entre las medias de las dos poblaciones es:

$$IC = (\bar{X}_1 - \bar{X}_2) \pm Z_{\alpha/2} \cdot \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Donde:

- $\bar{X}_1$  y  $\bar{X}_2$  son las medias de las muestras 1 y 2.
- $s_1^2$  y  $s_2^2$  son las varianzas de las muestras.
- $n_1$  y  $n_2$  son los tamaños de las muestras.
- $Z_{\alpha/2}$  es el valor crítico de la distribución normal estándar para un nivel de confianza dado (por ejemplo, 1.96 para un intervalo de confianza del 95%).

### 2.6.3 Ejemplo de Aplicación

Supongamos que queremos comparar el consumo calórico diario promedio entre hombres y mujeres de una ciudad. Tomamos dos muestras aleatorias independientes de hombres y mujeres:

- $n_1 = 100$  hombres, con una media muestral de  $\bar{X}_1 = 2500$  calorías y una desviación estándar  $s_1 = 300$  calorías.
- $n_2 = 120$  mujeres, con una media muestral de  $\bar{X}_2 = 2300$  calorías y una desviación estándar  $s_2 = 250$  calorías.

Queremos calcular un intervalo de confianza del 95% para la diferencia entre las medias de las dos poblaciones.

Cálculo

- Diferencia entre las medias muestrales:

$$\bar{X}_1 - \bar{X}_2 = 2500 - 2300 = 200$$

- Nivel de confianza: Para un intervalo de confianza del 95%, el valor crítico es  $Z_{\alpha/2} = 1,96$ .
- Error estándar de la diferencia entre las medias:

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = \sqrt{\frac{300^2}{100} + \frac{250^2}{120}}$$

$$SE = \sqrt{\frac{90000}{100} + \frac{62500}{120}} = \sqrt{900 + 520,83} = \sqrt{1420,83} = 37,7$$

- Margen de error:

$$ME = Z_{\alpha/2} \times SE = 1,96 \times 37,7 = 73,9$$

- Intervalo de confianza:

$$IC = (200 \pm 73,9)$$

$$IC = (126,1, 273,9)$$

Por lo tanto, podemos decir con un 95% de confianza que la diferencia entre el consumo calórico promedio de hombres y mujeres está entre 126.1 y 273.9 calorías. Como el intervalo no incluye el valor 0, podemos concluir que hay una diferencia significativa en el consumo calórico promedio entre hombres y mujeres.

El intervalo de confianza para la diferencia entre medias es una herramienta poderosa para realizar comparaciones entre dos poblaciones. Proporciona un rango de valores plausibles para la diferencia entre las medias y permite evaluar si existe una diferencia significativa entre los grupos, todo con un nivel de confianza específico.

Este método es ampliamente utilizado en investigaciones científicas, estudios de mercado y análisis de datos, y es fundamental en el proceso de toma de decisiones basada en datos.

## 2.7 Muestreo Estratificado

### 2.7.1 Definición

El muestreo estratificado es una técnica de muestreo en la que la población se divide en subgrupos homogéneos llamados “estratos”, en función de una o más características relevantes (por ejemplo, edad, género, nivel educativo). Luego, se seleccionan muestras aleatorias de cada uno de los estratos. Esta técnica asegura que todos los estratos estén representados en la muestra y mejora la precisión de las estimaciones, especialmente cuando hay diferencias significativas entre los estratos en relación con la variable de interés.

### 2.7.2 Tipos de Muestreo Estratificado

- Muestreo Estratificado Proporcional: El tamaño de la muestra seleccionada en cada estrato es proporcional al tamaño del estrato en la población total.
- Muestreo Estratificado No Proporcional: La muestra de cada estrato no tiene por qué ser proporcional al tamaño del estrato en la población.

### 2.7.3 Ventajas del Muestreo Estratificado

- Mayor precisión en las estimaciones poblacionales.
- Representación adecuada de subgrupos importantes.
- Mejora la comparación entre subgrupos.

### 2.7.4 Desventajas del Muestreo Estratificado

- Mayor complejidad logística.
- Difícil de aplicar cuando los estratos no están claramente definidos.

### 2.7.5 Pasos para Implementar el Muestreo Estratificado

1. Dividir la población en estratos homogéneos.
2. Determinar el tamaño de la muestra en cada estrato.
3. Seleccionar una muestra aleatoria dentro de cada estrato.
4. Combinar las muestras de los diferentes estratos.

### 2.7.6 Fórmula para el Tamaño de Muestra en Muestreo Estratificado Proporcional

$$n_h = \frac{N_h}{N} \cdot n$$

Donde:

- $N_h$  es el tamaño del estrato  $h$ ,
- $N$  es el tamaño total de la población,
- $n$  es el tamaño total de la muestra,
- $n_h$  es el tamaño de la muestra en el estrato  $h$ .

### 2.7.7 Ejemplo de Muestreo Estratificado

- Supongamos que queremos estudiar el nivel de actividad física entre los empleados de una empresa con 1,000 trabajadores. Dividimos a la población en tres estratos: producción, ventas y administración. Queremos una muestra de 200 empleados.
- Aplicamos la fórmula para calcular el tamaño de la muestra en cada estrato.

El muestreo estratificado es una técnica muy útil cuando se quiere asegurar que los diferentes subgrupos de una población estén representados en un estudio. Aunque su implementación es más compleja que el muestreo aleatorio simple, sus beneficios en términos de precisión y representatividad son evidentes.

## 2.8 Principio Aditivo, Multiplicativo y Arreglo Rectangular

### 2.8.1 Definición

Los principios aditivo y multiplicativo son reglas básicas de conteo utilizadas en combinatoria y probabilidad para calcular el número de formas en que pueden ocurrir diferentes eventos. El arreglo rectangular es una representación visual que se utiliza para organizar y contar las combinaciones posibles de varios conjuntos de elementos.

Estos principios son fundamentales para entender cómo se pueden combinar diferentes opciones y calcular el número total de resultados posibles en una situación dada.

### 2.8.2 Principio Aditivo

El principio aditivo se aplica cuando dos eventos son mutuamente excluyentes, es decir, no pueden ocurrir al mismo tiempo. El número total de formas en que cualquiera de los eventos

puede ocurrir es la suma de las formas en que puede ocurrir cada evento individualmente.

Si A y B son eventos mutuamente excluyentes, entonces:  $N(A \text{ o } B) = N(A) + N(B)$

Ejemplo del Principio Aditivo

Supongamos que un restaurante ofrece dos opciones de menú: desayuno o almuerzo.

- El menú de desayuno tiene 4 opciones.
- El menú de almuerzo tiene 6 opciones.

Si alguien puede elegir entre desayuno o almuerzo, el número total de opciones es:

$$N(\text{total}) = 4 + 6 = 10$$

Esto significa que hay 10 maneras diferentes de seleccionar una comida.

### 2.8.3 Principio Multiplicativo

El principio multiplicativo se utiliza cuando se quiere contar el número de formas en que dos o más eventos independientes pueden ocurrir en sucesión. Si hay  $n$  formas de realizar un primer evento y  $m$  formas de realizar un segundo evento, entonces el número total de formas en que ambos eventos pueden ocurrir es el producto de  $n$  y  $m$ .

Si A y B son eventos independientes, entonces:  $N(A \text{ y } B) = N(A) \times N(B)$

Ejemplo del Principio Multiplicativo

Supongamos que un menú ofrece 3 opciones de entrada y 4 opciones de plato principal. El número total de combinaciones posibles de entrada y plato principal es:

$$N(\text{total}) = 3 \times 4 = 12$$

Esto significa que hay 12 combinaciones posibles entre entrada y plato principal.

### 2.8.4 Arreglo Rectangular

Un arreglo rectangular es una representación gráfica que muestra todas las combinaciones posibles de dos o más conjuntos de elementos en un formato de tabla o cuadrícula. Este método es útil para visualizar cómo se combinan diferentes opciones y facilita el uso de los principios aditivo y multiplicativo.

Ejemplo de Arreglo Rectangular

Supongamos que en una tienda de ropa se pueden combinar 2 tipos de camisas (roja y azul) con 3 tipos de pantalones (negro, gris y blanco).

Camisa/Pantalón	Negro	Gris	Blanco
Roja	1	1	1
Azul	1	1	1

El número total de combinaciones posibles de camisas y pantalones es:

$$N(\text{total}) = 2 \times 3 = 6$$

El arreglo rectangular nos muestra las 6 combinaciones posibles de camisas y pantalones.

### 2.8.5 Relación entre Principio Aditivo, Multiplicativo y Arreglo Rectangular

Los principios aditivo y multiplicativo se aplican en diferentes contextos. El principio aditivo se utiliza cuando los eventos son mutuamente excluyentes y solo puede ocurrir uno de ellos, mientras que el principio multiplicativo se utiliza cuando se combinan diferentes eventos. El arreglo rectangular es una herramienta visual que ayuda a aplicar el principio multiplicativo al mostrar todas las combinaciones posibles de diferentes conjuntos de opciones.

### 2.8.6 Ejemplo Combinado de Principio Aditivo y Multiplicativo

Problema: Supongamos que en una tienda se ofrecen dos promociones:

- La primera promoción ofrece una camisa gratis si se compran pantalones, y hay 4 tipos de pantalones disponibles.
- La segunda promoción ofrece un par de zapatos gratis si se compran gafas, y hay 5 tipos de gafas disponibles.

Solución usando el principio aditivo y multiplicativo:

- Promoción 1:  $1 \times 4 = 4$  combinaciones de pantalones y camisas.
- Promoción 2:  $1 \times 5 = 5$  combinaciones de gafas y zapatos.
- Combinación total usando el principio aditivo:  $4 + 5 = 9$ .

El principio aditivo y el principio multiplicativo son herramientas fundamentales para contar las diferentes combinaciones posibles de eventos o elecciones, dependiendo de si los eventos son excluyentes o independientes. El arreglo rectangular proporciona una forma visual de representar combinaciones y aplicar estos principios de manera más clara.

## 2.9 Diagrama de Árbol y Principio Multiplicativo

### 2.9.1 Definición del Principio Multiplicativo

El principio multiplicativo es una de las reglas básicas de conteo utilizada en combinatoria y probabilidad para calcular el número total de formas en que varios eventos independientes pueden ocurrir en sucesión. Si un evento  $A$  puede ocurrir de  $n$  formas y un evento  $B$  puede ocurrir de  $m$  formas, entonces el número total de formas en que ambos eventos pueden ocurrir es el producto de  $n \times m$ .

$$N(\text{total}) = N(A) \times N(B) \times N(C) \times \dots$$

Ejemplo del Principio Multiplicativo

Supongamos que en un restaurante hay 3 opciones de entrantes, 4 opciones de platos principales y 2 opciones de postres.

$$N(\text{total}) = 3 \times 4 \times 2 = 24$$

Esto significa que hay 24 combinaciones diferentes de menú posibles.

### 2.9.2 Definición del Diagrama de Árbol

El diagrama de árbol es una representación gráfica utilizada para visualizar todas las combinaciones posibles de diferentes eventos o elecciones secuenciales. Es una herramienta que

nos permite aplicar el principio multiplicativo de manera visual, ya que muestra todas las ramas posibles que parten de cada elección.

Cada rama del árbol representa una elección o un resultado, y al final del diagrama podemos contar cuántos caminos posibles existen para llegar a una determinada combinación de elecciones.

### 2.9.3 Construcción de un Diagrama de Árbol

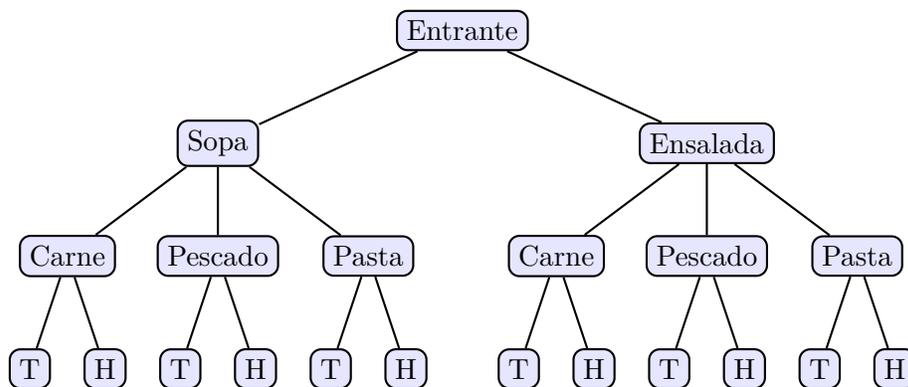
1. Comienza con el primer evento: Dibuja una rama para cada posible resultado del primer evento.
2. Añade ramas para el segundo evento: Desde el final de cada rama del primer evento, dibuja ramas adicionales para representar todas las opciones del segundo evento.
3. Continúa con eventos adicionales: Repite este proceso para cada evento, añadiendo ramas para cada elección disponible.

Ejemplo de Diagrama de Árbol

Imagina que vas a un restaurante y tienes las siguientes opciones:

- Entrante: sopa o ensalada.
- Plato principal: carne, pescado o pasta.
- Postre: tarta o helado.

El diagrama de árbol correspondiente sería:



Donde: T es igual a tarta, H es igual a helado

Cada camino desde la raíz (el entrante) hasta las hojas (el postre) representa una combinación única de entrante, plato principal y postre. Contando las hojas del diagrama de árbol, vemos que hay 12 combinaciones posibles:

$$N(\text{total}) = 2 \times 3 \times 2 = 12$$

### 2.9.4 Aplicaciones del Diagrama de Árbol

- Visualización de combinaciones: Los diagramas de árbol son útiles para organizar visualmente las combinaciones posibles en situaciones donde existen múltiples elecciones o eventos.

- Resolución de problemas de probabilidad: Se utilizan comúnmente en problemas de probabilidad para ilustrar los diferentes resultados de eventos secuenciales y calcular las probabilidades de eventos complejos.
- Planificación de decisiones: Son útiles para visualizar posibles escenarios y elecciones en proyectos o actividades.

### 2.9.5 Relación entre Diagrama de Árbol y Principio Multiplicativo

El principio multiplicativo y el diagrama de árbol están estrechamente relacionados, ya que el principio multiplicativo calcula cuántas combinaciones posibles hay de varios eventos, mientras que el diagrama de árbol proporciona una representación visual de esas combinaciones.

### 2.9.6 Ejemplo de Aplicación Combinada

Problema: Supongamos que estamos organizando una feria de comida y los clientes pueden elegir entre:

- 3 tipos de bebidas: refresco, agua, jugo.
- 2 tipos de bocadillos: sándwich o pizza.
- 2 tipos de postres: galletas o pastel.

Podemos usar el principio multiplicativo para calcular el número total de combinaciones posibles:

$$N(\text{total}) = 3 \times 2 \times 2 = 12$$

Además, podemos usar un diagrama de árbol para visualizar todas las combinaciones posibles de bebida, bocadillo y postre.

### 2.9.7 Ventajas de Usar Diagramas de Árbol

- Claridad visual: Permite visualizar todas las combinaciones posibles de eventos o decisiones secuenciales.
- Simplicidad: Hace que problemas complejos con múltiples elecciones sean más fáciles de entender.
- Organización: Ayuda a organizar las opciones de manera estructurada.

### 2.9.8 Desventajas de Usar Diagramas de Árbol

- Complejidad en casos grandes: Si hay muchas opciones, los diagramas de árbol pueden volverse grandes y difíciles de manejar visualmente.
- Menos eficiente para contar: El principio multiplicativo es más eficiente para contar el número total de combinaciones sin necesidad de dibujar todo el árbol.

### 2.9.9 Aplicaciones Prácticas del Principio Multiplicativo y los Diagramas de Árbol

- En educación: Para enseñar combinatoria, probabilidad y procesos de toma de decisiones.
- En negocios: Para planificar proyectos complejos que implican varias etapas o decisiones.

- En investigación: Para organizar y contar las combinaciones posibles de factores.
- En programación: Para representar estructuras de datos y algoritmos.

El principio multiplicativo es una herramienta fundamental en combinatoria para contar el número total de combinaciones posibles cuando ocurren varios eventos independientes en sucesión. Los diagramas de árbol proporcionan una forma visual de aplicar este principio y de representar las combinaciones posibles de una manera clara y estructurada.

## 2.10 Permutaciones

### 2.10.1 Definición

Las permutaciones son arreglos o secuencias de elementos en un orden específico. Se utilizan cuando el orden de los elementos es importante. Por ejemplo, en una carrera, los diferentes lugares ocupados por los corredores representan diferentes permutaciones de los corredores.

### 2.10.2 Fórmula para Permutaciones

Para calcular el número de permutaciones de  $n$  elementos tomados de  $r$  en  $r$ , se utiliza la siguiente fórmula:

$$P(n, r) = \frac{n!}{(n-r)!}$$

Donde:

- $P(n, r)$  es el número de permutaciones de  $n$  elementos tomados  $r$  a la vez.
- $n!$  (factorial de  $n$ ) es el producto de todos los números enteros desde 1 hasta  $n$ .
- $(n-r)!$  es el factorial de la diferencia entre el número total de elementos y el número de elementos seleccionados.

### 2.10.3 Ejemplo de Permutaciones

Supongamos que tenemos 4 libros diferentes (A, B, C, D) y queremos saber de cuántas maneras se pueden organizar 2 de esos libros en una estantería.

Aplicando la fórmula

- $n = 4$  (número total de libros)
- $r = 2$  (número de libros que se van a organizar)

$$P(4, 2) = \frac{4!}{(4-2)!} = \frac{4!}{2!} = \frac{4 \times 3 \times 2 \times 1}{(2 \times 1)} = 12$$

Listando las permutaciones

Las permutaciones son:

- AB
- AC
- AD
- BA

- BC
- BD
- CA
- CB
- CD
- DA
- DB
- DC

#### 2.10.4 Permutaciones de un Conjunto Completo

Cuando queremos calcular las permutaciones de todos los elementos de un conjunto, simplemente usamos el factorial del número total de elementos:

$$P(n) = n!$$

Por ejemplo, si tenemos 5 libros (A, B, C, D, E) y queremos saber de cuántas maneras se pueden organizar todos ellos:

$$P(5) = 5! = 5 \times 4 \times 3 \times 2 \times 1 = 120$$

#### 2.10.5 Permutaciones con Elementos Repetidos

Cuando algunos elementos en el conjunto son idénticos, se utiliza la siguiente fórmula:

$$P(n; n_1, n_2, \dots, n_k) = \frac{n!}{n_1! \times n_2! \times \dots \times n_k!}$$

Donde  $n$  es el número total de elementos y  $n_1, n_2, \dots, n_k$  son las cantidades de los elementos idénticos.

Ejemplo

Supongamos que tenemos las letras de la palabra “BANANA”, donde hay 6 letras en total y la letra “A” aparece 3 veces y la letra “N” aparece 2 veces. Entonces, el número de permutaciones es:

$$P(6; 3, 2, 1) = \frac{6!}{3! \times 2! \times 1!} = \frac{720}{6 \times 2 \times 1} = 60$$

#### 2.10.6 Aplicaciones de las Permutaciones

- Combinaciones de Ropa: Al elegir un atuendo donde el orden importa.
- Códigos y Contraseñas: En la generación de contraseñas o códigos PIN.
- Organización de Competencias: En competiciones deportivas.
- Juegos de Mesa: En juegos como el ajedrez.

### 2.10.7 Relación entre Permutaciones y Combinaciones

- Permutaciones: Se utilizan cuando el orden de los elementos es importante.
- Combinaciones: Se utilizan cuando el orden no es relevante.

Las permutaciones son una herramienta fundamental en combinatoria que nos permite contar de manera efectiva las diferentes formas en que se pueden organizar elementos. Entender cómo calcular y aplicar las permutaciones es esencial en diversas áreas como la probabilidad, la estadística y la planificación de eventos.

## 2.11 Combinaciones

### 2.11.1 Definición

Las combinaciones son selecciones de elementos de un conjunto en las que el orden de los elementos no importa. En otras palabras, al formar una combinación, los elementos elegidos son considerados como un grupo sin que se tenga en cuenta su disposición o secuencia.

### 2.11.2 Fórmula para Combinaciones

Para calcular el número de combinaciones de  $n$  elementos tomados de  $r$  en  $r$ , se utiliza la siguiente fórmula:

$$C(n, r) = \frac{n!}{r!(n-r)!}$$

Donde:

- $C(n, r)$  es el número de combinaciones de  $n$  elementos tomados  $r$  a la vez.
- $n!$  es el factorial de  $n$ .
- $r!$  es el factorial de  $r$ .
- $(n - r)!$  es el factorial de la diferencia entre el número total de elementos y el número de elementos seleccionados.

### 2.11.3 Ejemplo de Combinaciones

Supongamos que tenemos 5 frutas diferentes: manzana, plátano, naranja, uva y pera. Queremos saber de cuántas maneras podemos seleccionar 3 frutas de este grupo.

Aplicando la fórmula

- $n = 5$  (número total de frutas)
- $r = 3$  (número de frutas que se van a seleccionar)

$$C(5, 3) = \frac{5!}{3!(5-3)!} = \frac{5!}{3! \times 2!} = \frac{5 \times 4}{2 \times 1} = 10$$

Listando las combinaciones

Las combinaciones posibles de 3 frutas son:

- Manzana, plátano, naranja
- Manzana, plátano, uva

- Manzana, plátano, pera
- Manzana, naranja, uva
- Manzana, naranja, pera
- Manzana, uva, pera
- Plátano, naranja, uva
- Plátano, naranja, pera
- Plátano, uva, pera
- Naranja, uva, pera

Así, hay 10 maneras diferentes de seleccionar 3 frutas de un total de 5.

#### 2.11.4 Combinaciones de un Conjunto Completo

Cuando queremos calcular las combinaciones de todos los elementos de un conjunto, simplemente usamos la fórmula con  $r = n$ :

$$C(n, n) = 1$$

Esto significa que hay solo una forma de seleccionar todos los elementos del conjunto.

#### 2.11.5 Combinaciones con Elementos Repetidos

Cuando se permite que algunos elementos se repitan en la selección, se utiliza la siguiente fórmula:

$$C(n, r) = \frac{(n + r - 1)!}{r!(n - 1)!}$$

Donde:

- $n$  es el número de tipos de elementos,
- $r$  es el número total de elementos seleccionados.

Ejemplo de Combinaciones con Repetición

Supongamos que tienes 3 tipos de helados (vainilla, chocolate y fresa) y deseas seleccionar 4 bolas de helado. Permitiendo repeticiones, el número de combinaciones es:

$$C(3, 4) = \frac{(3 + 4 - 1)!}{4!(3 - 1)!} = \frac{6!}{4! \times 2!} = \frac{720}{24 \times 2} = 15$$

#### 2.11.6 Aplicaciones de las Combinaciones

- Selección de Equipos: Cuando se forma un equipo a partir de un grupo de personas.
- Elección de Menús: Al elegir un conjunto de platos de un menú.
- Análisis de Datos: En estudios estadísticos donde se requiere seleccionar una muestra.
- Loterías y Juegos de Azar: Cuando se seleccionan números o combinaciones de elementos.

### 2.11.7 Relación entre Combinaciones y Permutaciones

- Combinaciones: Se utilizan cuando el orden de los elementos no importa.
- Permutaciones: Se utilizan cuando el orden de los elementos es importante.

### 2.11.8 Ejemplo de Aplicación Combinada

Problema: En una encuesta, un investigador desea elegir 4 de las 10 preguntas para una prueba. ¿Cuántas combinaciones de preguntas son posibles?

Solución usando la fórmula:

$$C(10,4) = \frac{10!}{4!(10-4)!} = \frac{10!}{4! \times 6!} = \frac{10 \times 9 \times 8 \times 7}{4 \times 3 \times 2 \times 1} = 210$$

Las combinaciones son un concepto clave en combinatoria y probabilidades, que permiten calcular el número de maneras en que se pueden seleccionar elementos de un conjunto sin tener en cuenta el orden. Entender cómo calcular y aplicar las combinaciones es esencial en diversas disciplinas, desde la estadística hasta la planificación de eventos.



## 3 — Asociación Estadística entre Variables

### 3.1 Asociación Estadística entre Variables

#### 3.1.1 ¿Qué es una asociación entre variables?

Cuando hablamos de “asociación” en estadística, nos referimos a la relación que puede existir entre dos o más variables. Esto significa que si una variable cambia, es probable que la otra también lo haga de alguna manera predecible.

Ejemplo de la vida real:

Imagina que estás investigando la relación entre las horas que una persona estudia y la calificación que obtiene en un examen. Si las personas que estudian más tienden a obtener mejores calificaciones, podemos decir que hay una asociación positiva entre el tiempo de estudio y las calificaciones.

#### 3.1.2 Diferencia entre asociación y causalidad

Es crucial entender que una asociación no implica necesariamente que una variable cause el cambio en la otra. A veces, dos variables pueden estar relacionadas simplemente porque están influenciadas por otra variable que no estamos observando.

Ejemplo de correlación sin causalidad:

Podemos observar una relación entre el aumento de las ventas de helado y el aumento de la temperatura en verano. A más calor, más gente compra helado. Sin embargo, esto no significa que el consumo de helado cause que la temperatura suba. En este caso, hay una asociación entre el consumo de helado y la temperatura, pero no una causalidad.

### 3.1.3 Conceptos clave

- Asociación positiva: Cuando una variable aumenta, la otra también lo hace. Ejemplo: Altura y peso. Generalmente, a mayor altura, mayor peso.
- Asociación negativa: Cuando una variable aumenta, la otra disminuye. Ejemplo: Cantidad de ejercicio y porcentaje de grasa corporal. A más ejercicio, menos grasa corporal.
- No asociación: No hay ninguna relación entre dos variables. Ejemplo: El color de los zapatos y las notas en un examen.

Ejemplo práctico de asociación:

Supongamos que tienes los siguientes datos de cinco estudiantes:

Estudiante	Horas de estudio	Calificación
A	2	60
B	4	70
C	6	75
D	8	85
E	10	90

Podemos ver que, en general, conforme las horas de estudio aumentan, las calificaciones también lo hacen. Esto sugiere una asociación positiva entre las horas de estudio y las calificaciones.

## 3.2 Midiendo la Asociación entre Dos Variables

Para entender mejor la asociación entre dos variables, los estadísticos utilizan medidas que nos permiten cuantificar qué tan fuerte o débil es esta relación.

### 3.2.1 Coeficiente de correlación de Pearson

El coeficiente de correlación de Pearson mide la relación lineal entre dos variables numéricas (cuantitativas). Este coeficiente se representa con la letra  $r$  y siempre está entre -1 y 1:

- $r = 1$  significa una correlación positiva perfecta.
- $r = -1$  significa una correlación negativa perfecta.
- $r = 0$  significa que no hay correlación lineal entre las variables.

La fórmula para calcular el coeficiente de Pearson es la siguiente:

$$r = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum(X_i - \bar{X})^2 \sum(Y_i - \bar{Y})^2}}$$

Donde:

- $X_i$  e  $Y_i$  son los valores de las variables  $X$  y  $Y$ ,
- $\bar{X}$  y  $\bar{Y}$  son las medias de las variables  $X$  y  $Y$ ,
- $r$  es el coeficiente de correlación.

Ejemplo práctico de cálculo:

Supongamos que tienes los siguientes datos:

Estudiante	$X$ : Horas de estudio	$Y$ : Calificación
$A$	2	60
$B$	4	70
$C$	6	75
$D$	8	85
$E$	10	90

1. Primero, calculamos la media de las horas de estudio y la media de las calificaciones:

$$\bar{X} = \frac{2+4+6+8+10}{5} = 6, \quad \bar{Y} = \frac{60+70+75+85+90}{5} = 76$$

2. Luego, calculamos las desviaciones de cada valor respecto a la media:

$$X_i - \bar{X} : \quad -4, -2, 0, 2, 4$$

$$Y_i - \bar{Y} : \quad -16, -6, -1, 9, 14$$

3. Multiplicamos las desviaciones y sumamos los productos:

$$\sum (X_i - \bar{X})(Y_i - \bar{Y}) = (-4)(-16) + (-2)(-6) + (0)(-1) + (2)(9) + (4)(14) = 150$$

4. Finalmente, calculamos el denominador de la fórmula de Pearson y obtenemos:

$$r = \frac{150}{\sqrt{40 \times 740}} = \frac{150}{172,02} \approx 0,87$$

Esto indica una fuerte correlación positiva entre las horas de estudio y las calificaciones.

### 3.2.2 Correlación de Spearman

El coeficiente de correlación de Spearman es similar al de Pearson, pero se basa en los rangos de los datos en lugar de los valores numéricos exactos. Es útil cuando los datos no tienen una relación lineal clara o cuando hay valores atípicos que podrían distorsionar la medición de Pearson.

Ejemplo práctico:

Si tenemos datos que no siguen una línea recta, como la relación entre el estrés y el rendimiento académico (donde niveles moderados de estrés pueden mejorar el rendimiento, pero niveles altos lo empeoran), usaríamos Spearman en lugar de Pearson.

### 3.2.3 Interpretación de las correlaciones

Cuando calculamos el coeficiente de correlación (sea Pearson o Spearman), podemos interpretarlo de la siguiente manera:

- 0.7 a 1.0 (o -0.7 a -1.0): Correlación fuerte.
- 0.4 a 0.7 (o -0.4 a -0.7): Correlación moderada.
- 0.1 a 0.4 (o -0.1 a -0.4): Correlación débil.
- 0: No hay correlación.

Actividad recomendada:

Busca un conjunto de datos donde puedas medir la relación entre dos variables numéricas. Por ejemplo, podrías observar la relación entre la cantidad de horas que duermes cada noche y tu nivel de energía al día siguiente. Usa una calculadora o software como Excel para calcular el coeficiente de correlación de Pearson y analiza qué tan fuerte es la relación entre las dos variables.

## 3.3 El caso de dos variables categóricas

### 3.3.1 ¿Qué son las variables categóricas?

Las variables categóricas son aquellas que toman valores en forma de categorías o etiquetas. No tienen un orden numérico natural, y cada categoría representa una clase distinta. Por ejemplo, el color de los ojos (azul, verde, marrón) o el estado civil (soltero, casado, divorciado) son variables categóricas.

Cuando tratamos con dos variables categóricas, el objetivo es analizar si existe alguna relación o asociación entre ellas. Para hacerlo, podemos usar tablas de contingencia y medidas específicas que nos permitan cuantificar dicha asociación.

### 3.3.2 Tablas de contingencia

Una de las herramientas más útiles para analizar la relación entre dos variables categóricas es una tabla de contingencia (o tabla cruzada). Esta tabla muestra las frecuencias absolutas de las combinaciones de categorías de las dos variables.

Ejemplo práctico:

Supongamos que queremos analizar la relación entre el género (masculino o femenino) y la preferencia por un tipo de película (acción o comedia). Recopilamos datos de 100 personas y organizamos la información en la siguiente tabla de contingencia:

	Acción	Comedia	Total
Masculino	30	20	50
Femenino	10	40	50
Total	40	60	100

En esta tabla, podemos observar cómo se distribuyen las preferencias por tipo de película según el género. Este tipo de tabla nos ayuda a visualizar las relaciones entre dos variables categóricas.

### 3.3.3 Medidas de asociación para dos variables categóricas

Una de las medidas más comunes para analizar la asociación entre dos variables categóricas es el coeficiente de contingencia o el estadístico chi-cuadrado ( $\chi^2$ ). Este método se utiliza para determinar si las diferencias observadas en las frecuencias de las categorías son significativas o si se deben al azar.

La fórmula para el estadístico  $\chi^2$  es:

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

Donde:

- $O_i$  son las frecuencias observadas,
- $E_i$  son las frecuencias esperadas bajo la hipótesis de independencia.

Las frecuencias esperadas  $E_i$  se calculan usando la siguiente fórmula:

$$E_i = \frac{\text{Total fila} \times \text{Total columna}}{\text{Total general}}$$

Ejemplo práctico del cálculo de  $\chi^2$ :

Tomando el ejemplo anterior, calculemos las frecuencias esperadas para la celda de hombres que prefieren acción ( $O_1 = 30$ ):

$$E_1 = \frac{50 \times 40}{100} = 20$$

Hacemos esto para cada celda y luego calculamos el valor de  $\chi^2$ .

Si el valor de  $\chi^2$  es suficientemente grande, podemos concluir que hay una asociación entre las dos variables categóricas.

### 3.3.4 Interpretación

Si el valor de  $\chi^2$  es significativo, podemos afirmar que existe una relación entre las dos variables categóricas. En nuestro ejemplo, esto significaría que las preferencias por películas de acción o comedia están relacionadas con el género de las personas.

## 3.4 El caso de una variable categórica y una continua

### 3.4.1 ¿Qué son las variables continuas?

Las variables continuas son aquellas que pueden tomar un número infinito de valores dentro de un rango. Ejemplos comunes incluyen la altura, el peso o el ingreso anual. En este contexto, queremos analizar la relación entre una variable categórica (como género, estado civil, etc.) y una variable continua (como el ingreso o la edad).

Cuadro 3.1: Tabla de valores críticos chi-cuadrado

Grados de Libertad (df)	$\alpha = 0,05$	$\alpha = 0,01$
1	3.841	6.635
2	5.991	9.210
3	7.815	11.345
4	9.488	13.277
5	11.070	15.086
6	12.592	16.812
7	14.067	18.475
8	15.507	20.090
9	16.919	21.666
10	18.307	23.209

### 3.4.2 Comparación de medias

Cuando se estudia la relación entre una variable categórica y una variable continua, la manera más común de analizar esta asociación es comparando las medias de la variable continua en las diferentes categorías de la variable categórica.

Ejemplo práctico:

Supongamos que estamos interesados en comparar los ingresos anuales (variable continua) entre hombres y mujeres (variable categórica). Recopilamos los siguientes datos de ingresos (en miles de dólares) para cada grupo:

Género	Ingreso anual
Masculino	50, 60, 55, 70, 65
Femenino	40, 45, 50, 55, 48

Para comparar estos dos grupos, calculamos la media de cada grupo.

$$\text{Media hombres} = \frac{50 + 60 + 55 + 70 + 65}{5} = 60$$

$$\text{Media mujeres} = \frac{40 + 45 + 50 + 55 + 48}{5} = 47,6$$

Vemos que, en promedio, los hombres en esta muestra ganan más que las mujeres.

### 3.4.3 Prueba t de Student

Una forma de determinar si la diferencia entre las medias de dos grupos es significativa es usar la prueba t de Student. Esta prueba compara las medias de dos grupos para ver si las diferencias observadas podrían haber ocurrido por azar.

La fórmula para la prueba t es:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Donde:

- $\bar{X}_1$  y  $\bar{X}_2$  son las medias de los dos grupos,
- $s_1^2$  y  $s_2^2$  son las varianzas de los dos grupos,
- $n_1$  y  $n_2$  son los tamaños de muestra de los dos grupos.

Si el valor de  $t$  es lo suficientemente grande (dependiendo del nivel de significancia y los grados de libertad), podemos concluir que la diferencia entre las medias es estadísticamente significativa.

Ejemplo práctico de la prueba t:

Supongamos que calculamos las varianzas para ambos grupos:

$$s_1^2 = 62,5, \quad s_2^2 = 22,3$$

Y aplicamos la fórmula de la prueba t para comparar los ingresos de hombres y mujeres:

$$t = \frac{60 - 47,6}{\sqrt{\frac{62,5}{5} + \frac{22,3}{5}}} = \frac{12,4}{\sqrt{12,5 + 4,46}} = \frac{12,4}{\sqrt{16,96}} = \frac{12,4}{4,12} \approx 3,01$$

Con este valor de  $t$ , podemos verificar en tablas de distribución t si la diferencia es significativa.

#### 3.4.4 Interpretación

Si el valor de  $t$  es significativo, podemos concluir que existe una diferencia estadísticamente significativa entre los dos grupos. En nuestro ejemplo, esto podría significar que los ingresos de hombres y mujeres son realmente diferentes y no solo por azar.

#### 3.4.5 Conclusión

Hemos visto cómo analizar la relación entre dos variables categóricas mediante tablas de contingencia y el estadístico chi-cuadrado, así como la relación entre una variable categórica y una continua a través de la comparación de medias y la prueba t de Student. Estas herramientas son fundamentales para entender las asociaciones en diferentes tipos de datos.

### 3.5 El caso de dos variables cuantitativas

#### 3.5.1 ¿Qué son las variables cuantitativas?

Las variables cuantitativas son aquellas que se expresan en forma numérica y que pueden ser medidas o contadas. Ejemplos comunes de variables cuantitativas son la altura, el peso, el ingreso, entre otros. Estas variables permiten realizar operaciones aritméticas, como sumas y promedios.

Cuando analizamos la relación entre dos variables cuantitativas, generalmente queremos saber si existe una asociación entre ellas, es decir, si los cambios en una variable están relacionados con los cambios en la otra.

### 3.5.2 Diagrama de dispersión

Para visualizar la relación entre dos variables cuantitativas, utilizamos un diagrama de dispersión. Este gráfico muestra los pares de valores de las dos variables, colocando una variable en el eje  $X$  y la otra en el eje  $Y$ .

Ejemplo:

Supongamos que queremos analizar la relación entre el número de horas estudiadas y las calificaciones obtenidas en un examen. Los datos son los siguientes:

Estudiante	Horas de estudio (X)	Calificación (Y)
<i>A</i>	2	60
<i>B</i>	4	70
<i>C</i>	6	75
<i>D</i>	8	85
<i>E</i>	10	90

Podemos graficar estos datos en un diagrama de dispersión para observar la relación entre las dos variables. Si los puntos en el gráfico tienden a alinearse de manera ascendente, entonces existe una asociación positiva entre las horas de estudio y las calificaciones.

### 3.5.3 Coeficiente de correlación

Para cuantificar la fuerza y la dirección de la relación entre dos variables cuantitativas, utilizamos el coeficiente de correlación de Pearson. Este coeficiente, denotado por  $r$ , varía entre  $-1$  y  $1$ . Su valor indica lo siguiente:

- $r = 1$ : Correlación positiva perfecta.
- $r = -1$ : Correlación negativa perfecta.
- $r = 0$ : No hay correlación.

La fórmula para calcular  $r$  es:

$$r = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum(X_i - \bar{X})^2 \sum(Y_i - \bar{Y})^2}}$$

Ejemplo de cálculo:

Para los datos anteriores (horas de estudio y calificaciones), calculamos:

1.  $\bar{X} = 6$ ,  $\bar{Y} = 76$  2. Las desviaciones  $X_i - \bar{X}$  y  $Y_i - \bar{Y}$ :

$$X_i - \bar{X} = -4, -2, 0, 2, 4$$

$$Y_i - \bar{Y} = -16, -6, -1, 9, 14$$

3. Luego, calculamos el numerador y el denominador para obtener:

$$r \approx 0,87$$

Esto indica una fuerte correlación positiva entre las horas de estudio y las calificaciones.

### 3.6 El modelo de regresión lineal

#### 3.6.1 ¿Qué es la regresión lineal?

El modelo de regresión lineal es una herramienta estadística que nos permite modelar la relación entre dos variables cuantitativas. En su forma más simple, se utiliza para predecir el valor de una variable dependiente ( $Y$ ) en función del valor de una variable independiente ( $X$ ), bajo el supuesto de que la relación entre ellas es lineal.

La forma más simple de regresión lineal es la regresión lineal simple, que utiliza una sola variable independiente para predecir el valor de una variable dependiente. El modelo de regresión lineal simple se expresa mediante la ecuación de una recta:

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Donde:

- $Y$  es la variable dependiente.
- $X$  es la variable independiente.
- $\beta_0$  es la ordenada al origen (el valor de  $Y$  cuando  $X = 0$ ).
- $\beta_1$  es la pendiente de la recta (cuánto cambia  $Y$  por cada unidad que cambia  $X$ ).
- $\varepsilon$  es el término de error, que representa las variaciones que no son explicadas por el modelo.

#### 3.6.2 Estimación de los parámetros: método de los mínimos cuadrados

El método de los mínimos cuadrados es la técnica más utilizada para ajustar la recta de regresión a los datos observados. Este método encuentra los valores de  $\beta_0$  y  $\beta_1$  que minimizan la suma de los cuadrados de los residuales, es decir, la diferencia entre los valores observados y los valores predichos.

La fórmula para estimar la pendiente  $\beta_1$  es:

$$\beta_1 = \frac{\sum(X_i - \bar{X})(Y_i - \bar{Y})}{\sum(X_i - \bar{X})^2}$$

Y la fórmula para estimar la ordenada al origen  $\beta_0$  es:

$$\beta_0 = \bar{Y} - \beta_1 \bar{X}$$

Donde:

- $X_i$  e  $Y_i$  son los valores de las variables independientes y dependientes, respectivamente.
- $\bar{X}$  y  $\bar{Y}$  son las medias de  $X$  y  $Y$ , respectivamente.

Ejemplo de cálculo:

Para los datos de horas de estudio y calificaciones:

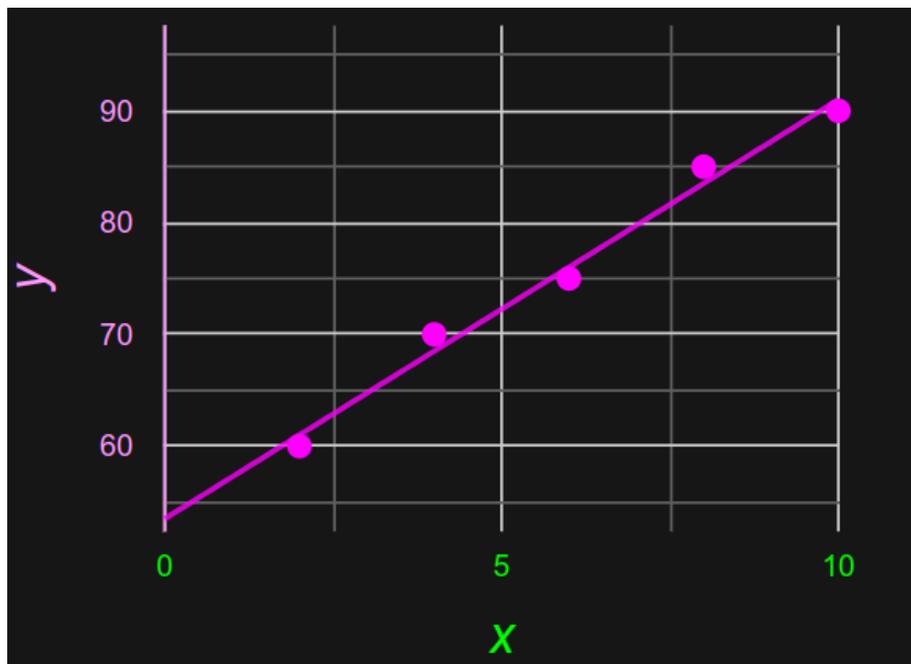
1. Hemos calculado que  $\bar{X} = 6$  y  $\bar{Y} = 76$ . 2. Utilizando los mismos datos de las desviaciones, calculamos:

$$\beta_1 = \frac{150}{40} = 3,75$$

$$\beta_0 = 76 - 3,75 \times 6 = 53,5$$

Por lo tanto, la ecuación de la recta de regresión es:

$$Y = 53,5 + 3,75X$$



### 3.6.3 Interpretación del modelo

La ecuación de la recta de regresión  $\hat{Y} = 53,5 + 3,75X$  nos dice que, en promedio:

- Por cada hora adicional de estudio, la calificación aumenta en 3.75 puntos.
- Si un estudiante no estudia ( $X = 0$ ), se espera que su calificación sea de 53.5 puntos.

### 3.6.4 Evaluación del modelo: coeficiente de determinación $R^2$

Una vez que se ha ajustado el modelo, es importante evaluar qué tan bien predice los valores observados. Para ello, podemos calcular el coeficiente de determinación  $R^2$ , que nos indica la proporción de la variabilidad en  $Y$  que es explicada por  $X$ . El valor de  $R^2$  varía entre 0 y 1, donde:

- $R^2 = 1$  indica un ajuste perfecto.

- $R^2 = 0$  indica que el modelo no explica la variabilidad en los datos.

La fórmula de  $R^2$  es:

$$R^2 = 1 - \frac{\sum(Y_i - \hat{Y}_i)^2}{\sum(Y_i - \bar{Y})^2}$$

Donde  $\hat{Y}_i$  son los valores predichos por el modelo.

Si calculamos  $R^2$  para los datos anteriores, supongamos que obtenemos  $R^2 = 0,9868$ , esto indica que el 98.68% de la variabilidad en las calificaciones es explicada por el número de horas de estudio, lo cual es un ajuste excelente.

### 3.6.5 Concepto de error o residual

El error o residual ( $\varepsilon$ ) es la diferencia entre el valor observado de  $Y$  y el valor predicho por el modelo. Se expresa como:

$$\varepsilon = Y - \hat{Y}$$

Donde  $\hat{Y}$  es el valor de  $Y$  predicho por la ecuación de regresión.

Ejemplo:

Supongamos que estamos interesados en modelar la relación entre la cantidad de horas de estudio ( $X$ ) y la calificación obtenida en un examen ( $Y$ ). Si los datos de un estudiante son  $X = 4$  horas y  $Y = 75$  puntos, pero el modelo de regresión predice una calificación de 78 puntos, entonces el error o residual es:

$$\varepsilon = 75 - 78 = -3$$

El valor negativo del residual indica que la predicción fue mayor que el valor observado.

### 3.6.6 Propiedades importantes del análisis de regresión lineal

Algunas propiedades clave del análisis de regresión lineal son:

- Los residuales tienen una media de cero, lo que significa que los errores positivos y negativos se compensan.
- La relación entre las variables  $X$  y  $Y$  es lineal.
- Se asume que los errores siguen una distribución normal y tienen varianza constante (homocedasticidad).

## 3.7 Bondad de ajuste del modelo de regresión

### 3.7.1 ¿Qué es la bondad de ajuste?

La bondad de ajuste de un modelo de regresión nos indica qué tan bien el modelo predice los valores observados de la variable dependiente  $Y$  a partir de los valores de la variable

independiente  $X$ . En otras palabras, mide qué tan bien los valores predichos por el modelo se ajustan a los datos reales.

La bondad de ajuste se evalúa mediante diversas métricas que cuantifican el grado en que los valores predichos por el modelo se aproximan a los valores observados. Las métricas más comunes para evaluar la bondad de ajuste de un modelo de regresión son:

- El coeficiente de determinación  $R^2$ .
- El error estándar de los residuos.

### 3.7.2 Error estándar de los residuos

El error estándar de los residuos es otra medida que evalúa qué tan bien se ajusta el modelo a los datos. Este mide la variabilidad de los valores observados de  $Y$  con respecto a los valores predichos por el modelo. El error estándar de los residuos se calcula como:

$$\text{Error estándar} = \sqrt{\frac{SCE}{n-2}}$$

Donde:

- $SCE$  es la suma de los cuadrados de los residuos.
- $n$  es el número de observaciones.

Ejemplo:

Para los datos anteriores, con  $SCE = 7,5$  y  $n = 5$ :

$$\text{Error estándar} = \sqrt{\frac{7,5}{5-2}} = \sqrt{\frac{7,5}{3}} \approx 1,58$$

Un error estándar bajo indica que los valores predichos por el modelo están muy cerca de los valores observados.

### 3.7.3 Conclusión sobre la bondad de ajuste

En este ejemplo, hemos utilizado tres métricas para evaluar la bondad de ajuste del modelo:

- El coeficiente de determinación  $R^2 = 0,9868$ , lo que indica que el modelo explica un 98.68% de la variabilidad de los datos.
- El error estándar de los residuos es 1.58, lo que indica una baja dispersión de los datos alrededor de la línea de regresión.

Estos resultados indican que el modelo de regresión ajusta muy bien a los datos, y es útil para hacer predicciones sobre la variable dependiente  $Y$  (calificaciones) en función de la variable independiente  $X$  (horas de estudio).

## 3.8 Teoría de la probabilidad

La teoría de la probabilidad es una rama de las matemáticas que se ocupa de los fenómenos aleatorios. Su objetivo es cuantificar la incertidumbre asociada con eventos que no pueden predecirse de manera determinista, pero sobre los cuales podemos establecer expectativas basadas en la información disponible.

### 3.8.1 Definiciones Básicas

- Experimento aleatorio: Un experimento cuyo resultado no puede predecirse con certeza. Ejemplo: lanzar un dado.
- Espacio muestral ( $S$ ): El conjunto de todos los posibles resultados de un experimento aleatorio. Ejemplo: Al lanzar un dado,  $S = \{1, 2, 3, 4, 5, 6\}$ .
- Evento ( $A$ ): Cualquier subconjunto del espacio muestral. Ejemplo: Obtener un número par al lanzar un dado,  $A = \{2, 4, 6\}$ .

### 3.8.2 Probabilidad Clásica

Si un experimento tiene  $n$  resultados igualmente probables y un evento  $A$  ocurre en  $m$  de esos resultados, la probabilidad de que ocurra  $A$  es:

$$P(A) = \frac{m}{n}$$

Ejemplo:

Al lanzar un dado, la probabilidad de obtener un número par es:

$$P(\text{par}) = \frac{3}{6} = 0,5$$

### 3.8.3 Propiedades de la Probabilidad

- $0 \leq P(A) \leq 1$ .
- La probabilidad de un evento seguro (evento que siempre ocurre) es 1:  $P(S) = 1$ .
- La probabilidad de un evento imposible es 0:  $P(\emptyset) = 0$ .
- Si  $A$  y  $B$  son eventos mutuamente excluyentes, entonces:

$$P(A \cup B) = P(A) + P(B)$$

## 3.9 Modelos teóricos de distribución de probabilidad

Un modelo de distribución de probabilidad es una función que asigna probabilidades a los posibles resultados de un experimento aleatorio. Estos modelos pueden ser discretos o continuos.

### 3.9.1 Distribuciones Discretas

Una distribución de probabilidad discreta es aquella en la que las variables aleatorias toman valores finitos o contables. Ejemplos: distribución binomial, distribución de Poisson.

### 3.9.2 Distribuciones Continuas

Una distribución de probabilidad continua es aquella en la que las variables aleatorias pueden tomar cualquier valor dentro de un rango continuo. Ejemplo: distribución normal.

## 3.10 La distribución binomial

La distribución binomial es un modelo discreto que describe el número de éxitos en una secuencia de  $n$  ensayos independientes, cada uno con una probabilidad de éxito  $p$ .

### 3.10.1 Función de probabilidad

La probabilidad de obtener exactamente  $k$  éxitos en  $n$  ensayos se calcula con la fórmula:

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

Donde:

- $n$  es el número de ensayos.
- $k$  es el número de éxitos deseados.
- $p$  es la probabilidad de éxito en cada ensayo.

Ejemplo:

Si lanzamos una moneda justa (probabilidad de cara  $p = 0,5$ ) 3 veces, la probabilidad de obtener exactamente 2 caras es:

$$P(X = 2) = \binom{3}{2} (0,5)^2 (0,5)^1 = 3 \times 0,25 \times 0,5 = 0,375$$

## 3.11 La distribución o curva normal

La distribución normal es una distribución continua que tiene una forma de campana simétrica. Es uno de los modelos más importantes en la probabilidad y estadística debido a su presencia en una gran cantidad de fenómenos naturales y sociales.

### 3.11.1 Función de densidad de probabilidad

La función de densidad de una variable aleatoria  $X$  que sigue una distribución normal con media  $\mu$  y desviación estándar  $\sigma$  está dada por:

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

Propiedades de la distribución normal:

- La curva es simétrica respecto a la media  $\mu$ .
- Aproximadamente el 68% de los valores se encuentran dentro de un intervalo de  $\pm 1$  desviación estándar de la media.
- Aproximadamente el 95% de los valores se encuentran dentro de un intervalo de  $\pm 2$  desviaciones estándar de la media.

### 3.11.2 Ejemplo

Si la altura de los estudiantes en una clase sigue una distribución normal con una media de 170 cm y una desviación estándar de 10 cm, podemos calcular la probabilidad de que un estudiante tenga una altura entre 160 cm y 180 cm. Usamos la tabla de la distribución normal estandarizada para obtener esta probabilidad.

Primero, estandarizamos los valores:

$$z = \frac{X - \mu}{\sigma}$$

Para  $X = 160$ :

$$z_1 = \frac{160 - 170}{10} = -1$$

Para  $X = 180$ :

$$z_2 = \frac{180 - 170}{10} = 1$$

De la tabla de la distribución normal, la probabilidad de estar entre  $-1$  y  $1$  es aproximadamente  $0.6826$ , lo que significa que el  $68.26\%$  de los estudiantes tienen una altura entre  $160$  cm y  $180$  cm.

### 3.12 La selección de la muestra

#### 3.12.1 ¿Qué es la selección de la muestra?

La selección de la muestra es el proceso mediante el cual se eligen individuos o unidades representativas de una población para realizar un estudio estadístico. La idea es que la muestra proporcione información suficiente sobre la población sin la necesidad de estudiar a toda la población.

#### 3.12.2 Tipos de muestreo

- Muestreo aleatorio simple: Cada elemento de la población tiene la misma probabilidad de ser seleccionado.
- Muestreo estratificado: La población se divide en grupos (estratos) y se toma una muestra aleatoria de cada estrato.
- Muestreo sistemático: Se seleccionan elementos de la población a intervalos regulares.

#### 3.12.3 Tamaño de la muestra

El tamaño de la muestra depende de varios factores, como el nivel de confianza deseado y el margen de error aceptable. Una fórmula común para calcular el tamaño de la muestra para estimar una proporción es:

$$n = \frac{Z^2 p(1-p)}{E^2}$$

Donde:

- $Z$  es el valor crítico asociado con el nivel de confianza.
- $p$  es la proporción esperada de la población que tiene la característica de interés.
- $E$  es el margen de error.

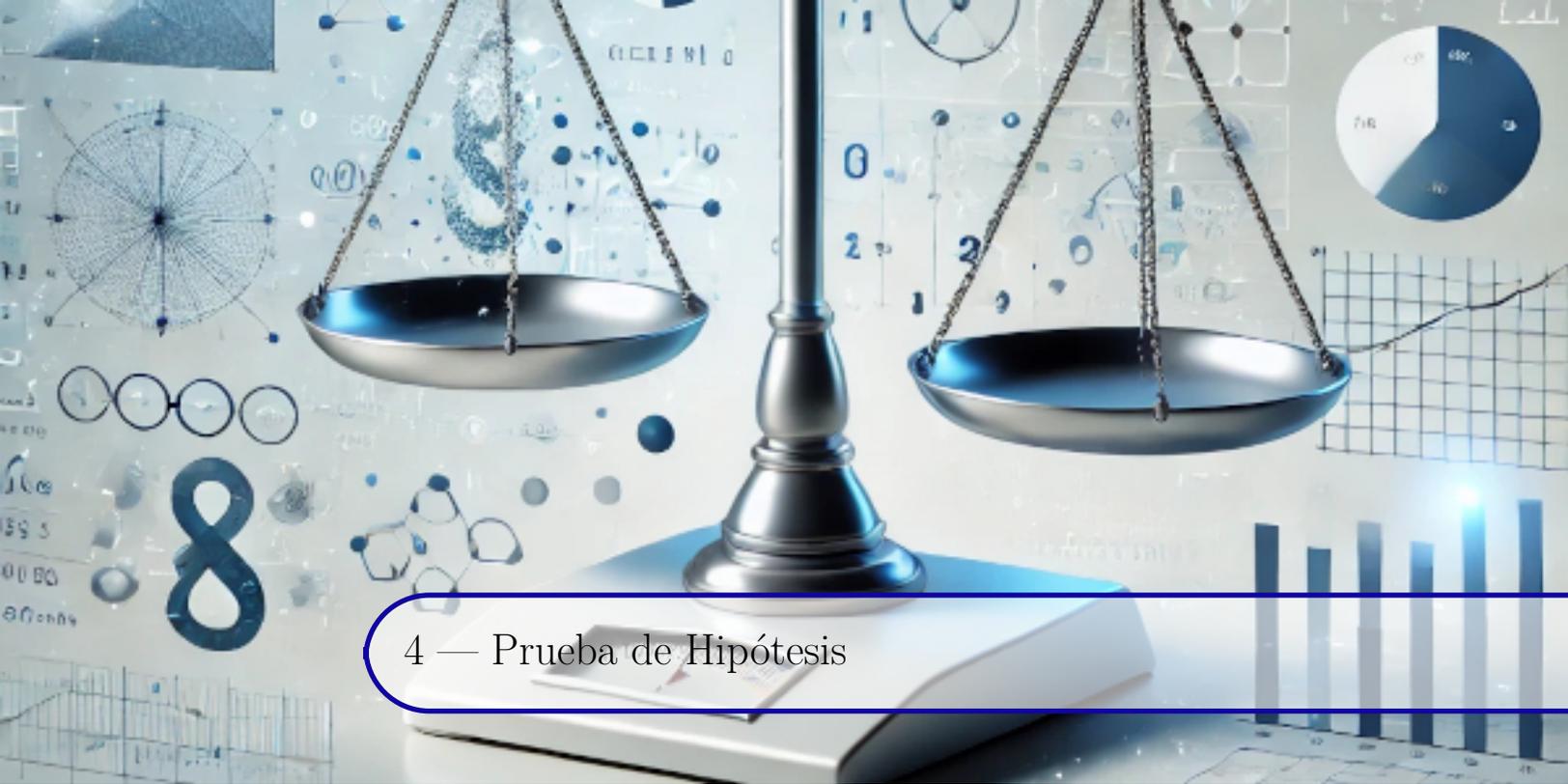
Ejemplo:

Si se quiere estimar la proporción de estudiantes que aprueban un examen con un nivel de confianza del  $95\%$  ( $Z = 1,96$ ) y un margen de error del  $5\%$  ( $E = 0,05$ ), y se espera que el  $50\%$  de los estudiantes aprueben ( $p = 0,5$ ), el tamaño de la muestra sería:

$$n = \frac{(1,96)^2 \times 0,5 \times (1 - 0,5)}{(0,05)^2} = \frac{3,8416 \times 0,25}{0,0025} = 384,16$$

Por lo tanto, se necesitaría una muestra de  $385$  estudiantes.





## 4 — Prueba de Hipótesis

### 4.1 Metodología para la prueba de hipótesis

Una hipótesis es una afirmación provisional que se somete a prueba para determinar su veracidad. La inferencia estadística establece un proceso para la prueba de hipótesis, que comienza con su enunciado formal y continúa con su contraste basado en la evidencia de los datos recolectados. Estos datos determinan si una hipótesis puede aceptarse como verdadera o debe rechazarse.

Este proceso se centra en los parámetros de la población (universo), y la hipótesis se formula con base en observaciones realizadas en una muestra representativa de esta población. Por ejemplo, una hipótesis como: “Estos pacientes tardan en promedio 25 días en recuperarse” implica que, en toda la población, el tiempo promedio de recuperación es de 25 días. El investigador debe contrastar esta afirmación comparando el valor hipotético (25 días) con los datos reales de una muestra. Si el promedio de la muestra resulta ser de 22 días, la estadística evaluará si esta diferencia permite aceptar o rechazar la hipótesis. El método estadístico será decisivo para resolver el dilema al evaluar la significación de la diferencia entre los valores.

¿Azar o no?

El método de pruebas de hipótesis consiste en determinar si una diferencia observada se debe al azar. Existen dos situaciones principales:

- a) Comparación entre un valor muestral y un valor poblacional o teórico.
- b) Comparación entre valores de dos o más muestras.

En el caso (a), se evalúa la diferencia entre un estadístico de la muestra y un parámetro de la población. En el caso (b), se evalúa la diferencia entre estadísticos de distintas muestras. Los valores comparados pueden ser promedios, porcentajes, etc. Nos enfocaremos en

promedios y porcentajes. La prueba estadística permite determinar la significación de estas diferencias con base en una hipótesis formulada estadísticamente.

Con la distribución de probabilidad adecuada, se calcula la probabilidad de observar dicha diferencia. Si esta probabilidad es baja, se considera que la diferencia es significativa.

## 4.2 Hipótesis nula y alternativa

Las hipótesis nula y alternativa son enunciados mutuamente excluyentes sobre una población. La prueba de hipótesis utiliza datos muestrales para decidir si se puede rechazar la hipótesis nula.

### 4.2.1 Hipótesis nula ( $H_0$ )

La hipótesis nula (denotada como  $H_0$ ) es una afirmación estadística que propone que no existe una diferencia significativa o un efecto particular en la población que se estudia, o bien, que cualquier diferencia observada es el resultado del azar o variabilidad natural de los datos.

En una prueba de hipótesis, la hipótesis nula es el punto de partida, y se asume como verdadera hasta que los datos muestrales proporcionen evidencia suficiente para rechazarla. Si se rechaza, indica que es probable que exista una diferencia real o un efecto en los datos.

Es la afirmación o suposición inicial que se pone a prueba. Se considera como la “hipótesis de no cambio” o “hipótesis de no efecto”.

Generalmente, la hipótesis nula sostiene que no existe una diferencia o efecto significativo entre los grupos o variables que se están comparando.

Ejemplo 1:

En una prueba de un medicamento, la hipótesis nula podría ser: “El medicamento no tiene efecto en la salud del paciente”.

Ejemplo 2:

Si queremos saber si el tiempo promedio de respuesta a una encuesta es de 10 minutos, la hipótesis nula sería:

$$H_0 : \mu = 10$$

donde  $\mu$  representa el tiempo promedio de respuesta en la población.

#### 4.2.2 Hipótesis alternativa ( $H_1$ )

La hipótesis alternativa (denotada como  $H_1$  o  $H_a$ ) es una afirmación opuesta a la hipótesis nula y propone que sí existe una diferencia significativa o un efecto en la población que se estudia. Esta hipótesis es lo que el investigador sospecha o espera probar como cierto. La hipótesis alternativa se considera verdadera si se obtiene suficiente evidencia estadística para rechazar la hipótesis nula.

Es la afirmación que se establece como opuesta a la hipótesis nula. Representa una posible conclusión que se puede aceptar si los datos proporcionan suficiente evidencia para rechazar la hipótesis nula.

En otras palabras, la hipótesis alternativa es la hipótesis que sugiere que sí hay un efecto o una diferencia.

Ejemplo: En el caso del medicamento, la hipótesis alternativa sería “El medicamento tiene un efecto positivo en la salud del paciente”.

Tipos de hipótesis alternativa

- Hipótesis bilateral: Sugiere que el parámetro de la población es diferente del valor especificado en la hipótesis nula, sin indicar una dirección específica. Se expresa como:

$$H_1 : \mu \neq \mu_0$$

donde  $\mu_0$  es el valor propuesto en la hipótesis nula.

- Hipótesis unilateral: Indica que el parámetro de la población es mayor o menor que el valor en la hipótesis nula, y tiene una dirección específica. Por ejemplo:
  - Hipótesis unilateral a la derecha:

$$H_1 : \mu > \mu_0$$

- Hipótesis unilateral a la izquierda:

$$H_1 : \mu < \mu_0$$

Ejemplo de hipótesis alternativa:

Siguiendo el ejemplo anterior, si se espera que el tiempo promedio de respuesta a una encuesta sea diferente de 10 minutos, la hipótesis alternativa sería:

$$H_1 : \mu \neq 10$$

Esta hipótesis implica que el tiempo promedio de respuesta es significativamente distinto de 10 minutos, en cualquiera de las dos direcciones (mayor o menor).

### 4.3 Error tipo I y tipo II

Al realizar pruebas de hipótesis, siempre existe la posibilidad de cometer un error. Existen dos tipos de errores:

#### 4.3.1 Error de tipo I

Un error de tipo I ocurre si se rechaza la hipótesis nula siendo verdadera. La probabilidad de este error es  $\alpha$ , el nivel de significancia de la prueba.

#### 4.3.2 Error de tipo II

Un error de tipo II ocurre si no se rechaza la hipótesis nula cuando es falsa. La probabilidad de este error es  $\beta$ , y su complemento ( $1 - \beta$ ) es la potencia de la prueba.

Decisión basada en la muestra	Verdad acerca de la población	
	H0 es verdadera	H0 es falsa
No rechazar H0	Decisión correcta (Probabilidad = $1 - \alpha$ )	Error tipo II - No rechazar H0 cuando es falsa (Probabilidad = $\beta$ )
Rechazar H0	Error tipo I - rechazar H0 cuando es verdadera (probabilidad = $\alpha$ )	Decisión correcta (probabilidad = $1 - \beta$ )

#### 4.3.3 Ejemplo

Supongamos que un investigador médico compara la efectividad de dos medicamentos:

- Hipótesis nula (H0):  $\mu_1 = \mu_2$  (los medicamentos tienen la misma eficacia).
- Hipótesis alternativa (H1):  $\mu_1 \neq \mu_2$  (los medicamentos tienen eficacias diferentes).

Un error de tipo I significa concluir que los medicamentos son diferentes cuando no lo son, mientras que un error de tipo II significa no detectar una diferencia real en su efectividad.

### 4.4 Pruebas de hipótesis

Al decidir entre dos hipótesis basadas en parámetros poblacionales, se debe definir el error aceptable (nivel de significancia). Las hipótesis se formulan así:

- H0 (hipótesis nula): Representa lo opuesto a lo que se espera confirmar, usando signos como  $\leq$  o  $\geq$ .
- H1 (hipótesis alternativa): Representa lo que se espera que sea cierto, usando signos  $>$  o  $<$ .

Se pueden cometer errores de tipo I o tipo II, controlados por los valores de  $\alpha$  y  $\beta$ . De ambos,  $\alpha$  es fundamental, ya que representa la probabilidad de rechazar la hipótesis nula erróneamente. Normalmente, se fija un valor de 5%.

### 4.5 Prueba de hipótesis Z para la media

Dentro de la inferencia estadística, la prueba Z permite estimar parámetros poblacionales a partir de una muestra, basándose en el teorema del límite central. Esta prueba calcula intervalos de confianza para determinar si el parámetro se encuentra dentro de un rango esperado, evaluando la validez de una aseveración sobre un parámetro poblacional. La Prueba de Hipótesis para una muestra es un método esencial para este análisis.

#### 4.5.1 Requisitos para realizar la Prueba Z

La prueba Z es una herramienta estadística que se utiliza para contrastar hipótesis sobre la media de una población. Los principales requisitos para aplicar la prueba Z son los siguientes:

1. Tamaño de la muestra ( $n$ ) grande: La prueba Z es más precisa cuando el tamaño de la muestra es mayor o igual a 30 ( $n \geq 30$ ). Esto se debe al Teorema del Límite Central, que asegura que la distribución de las medias muestrales tiende a una distribución normal.
2. Desviación estándar de la población conocida: Para utilizar la prueba Z, es necesario conocer la desviación estándar ( $\sigma$ ) de la población. Si no se conoce, se puede usar la desviación estándar muestral para muestras grandes como aproximación.
3. Muestreo aleatorio y representativo: La muestra debe ser seleccionada aleatoriamente y representar adecuadamente a la población de interés.
4. Distribución normal de la población: Si el tamaño de muestra es pequeño ( $n < 30$ ), se asume que la población sigue una distribución normal. Cuando la muestra es grande, la prueba Z puede aplicarse sin importar la distribución de la población.

#### 4.5.2 Cálculo de la Prueba Z y Ejemplo con Distribuciones Normales

La fórmula para calcular el valor Z en una prueba de hipótesis es:

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

donde:

- $\bar{X}$ : media muestral,
- $\mu$ : media poblacional bajo la hipótesis nula,
- $\sigma$ : desviación estándar de la población,
- $n$ : tamaño de la muestra.

El valor Z resultante permite comparar la media muestral con la media poblacional para determinar si la diferencia observada es estadísticamente significativa.

#### 4.5.3 Tabla de Valores Críticos Z para la Prueba de Hipótesis

En las pruebas de hipótesis, el valor crítico es el punto de corte que determina si se rechaza o no la hipótesis nula  $H_0$ . Este valor depende del nivel de significancia  $\alpha$  y de la distribución utilizada. En la prueba Z, los valores críticos corresponden a las áreas bajo la curva de la distribución normal estándar.

Existen dos tipos de pruebas que se utilizan con la distribución Z:

- Prueba de una cola (solo en una dirección)
- Prueba de dos colas (en ambas direcciones)

A continuación se presentan las tablas con los valores críticos Z para distintos niveles de significancia  $\alpha$  en ambos tipos de prueba.

#### 4.5.4 Tabla de Valores Críticos Z para la Prueba de Dos Colas

En una prueba de dos colas, el valor de  $\alpha$  se divide entre las dos colas de la distribución. Por ejemplo, para un nivel de significancia de  $\alpha = 0,05$ , el valor crítico Z es  $\pm 1,960$ .

$\alpha$	Valor Crítico Z (Dos Colas)
0.10	$\pm 1,645$
0.05	$\pm 1,960$
0.01	$\pm 2,576$
0.001	$\pm 3,291$

Cuadro 4.1: Valores Críticos Z para la Prueba de Dos Colas

#### 4.5.5 Tabla de Valores Críticos Z para la Prueba de Una Cola

En una prueba de una cola, todo el nivel de significancia  $\alpha$  se concentra en una sola cola de la distribución. Por ejemplo, para un nivel de significancia de  $\alpha = 0,05$ , el valor crítico Z es 1,645.

$\alpha$	Valor Crítico Z (Una Cola)
0.10	1.280
0.05	1.645
0.01	2.326
0.001	3.090

Cuadro 4.2: Valores Críticos Z para la Prueba de Una Cola

Los valores críticos Z son fundamentales para la toma de decisiones en las pruebas de hipótesis. Dependiendo de si se realiza una prueba de una cola o dos colas, y del nivel de significancia  $\alpha$ , estos valores determinan el umbral a partir del cual se rechaza la hipótesis nula.

#### 4.5.6 Caso 1: Prueba de dos colas

En una prueba de dos colas, el nivel de significancia  $\alpha$  se divide entre las dos colas de la distribución, por lo que cada cola tiene un área de  $\alpha/2$ .

Dado que  $\alpha = 0,05$ , el área de cada cola es:

$$\frac{\alpha}{2} = \frac{0,05}{2} = 0,025$$

Por lo tanto, el área acumulada en la cola superior es:

$$1 - 0,025 = 0,975$$

Buscamos el valor crítico correspondiente a una probabilidad acumulada de 0.975 en la tabla Z. El valor crítico Z para esta probabilidad es:

$$Z_{\alpha/2} = 1,960$$

El valor crítico en una prueba de dos colas es  $\pm 1,960$ . Como el valor Z calculado es 2.11, que es mayor que 1.960, rechazamos la hipótesis nula  $H_0$ .

#### 4.5.7 Caso 2: Prueba de una cola (a la derecha)

En una prueba de una cola, todo el nivel de significancia  $\alpha$  se concentra en una sola cola. Para  $\alpha = 0,05$ , buscamos el valor crítico para un área acumulada de:

$$1 - 0,05 = 0,95$$

El valor crítico Z correspondiente a una probabilidad acumulada de 0.95 es:

$$Z_{\alpha} = 1,645$$

En este caso, el valor crítico es  $Z = 1,645$ . Como el valor Z calculado es 2.11, que es mayor que 1.645, también rechazamos la hipótesis nula  $H_0$ .

#### 4.5.8 Ejemplo 1: Prueba de Hipótesis para la Media (Prueba Z Unilateral)

Supongamos que una empresa afirma que el tiempo promedio de recuperación con un medicamento es de 15 días, con una desviación estándar de 3 días en la población. Un investigador toma una muestra de 40 pacientes, encontrando un promedio de 16 días.

- Datos:  $\mu = 15$ ,  $\sigma = 3$ ,  $\bar{X} = 16$ ,  $n = 40$
- Nivel de significancia:  $\alpha = 0,05$

Cálculo del valor Z:

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{16 - 15}{\frac{3}{\sqrt{40}}} = \frac{1}{0,4743} \approx 2,11$$

#### 4.5.9 Prueba Z y Valor Crítico

En las pruebas de hipótesis, el valor crítico se utiliza para determinar si se debe rechazar la hipótesis nula  $H_0$  o no, dependiendo de la estadística calculada y el nivel de significancia  $\alpha$ . En este caso, consideramos una prueba Z con un valor Z calculado de 2.11 y un nivel de significancia de  $\alpha = 0,05$ .

Para una prueba Z con un valor Z calculado de 2.11 y un nivel de significancia  $\alpha = 0,05$ :

- En una prueba de dos colas, rechazamos la hipótesis nula  $H_0$  porque  $Z = 2,11 > 1,960$ .
- En una prueba de una cola (a la derecha), también rechazamos  $H_0$  porque  $Z = 2,11 > 1,645$ .

En ambos casos el cálculo está sugiriendo que el tiempo promedio de recuperación es mayor a 15 días.

#### 4.5.10 Ejemplo 2: Prueba de Hipótesis para la Media (Prueba Z Bilateral)

Un fabricante de bombillas asegura que la vida útil promedio es de 1200 horas con una desviación estándar de 100 horas. Un investigador toma una muestra de 50 bombillas, obteniendo una media de 1175 horas.

- Datos:  $\mu = 1200$ ,  $\sigma = 100$ ,  $\bar{X} = 1175$ ,  $n = 50$

- Nivel de significancia:  $\alpha = 0,05$  (valor crítico de  $Z$  es  $\pm 1.96$ ).

Cálculo del valor  $Z$ :

$$Z = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{1175 - 1200}{\frac{100}{\sqrt{50}}} = \frac{-25}{14,14} \approx -1,77$$

Interpretación: El valor  $Z = -1,77$  cae dentro del rango de aceptación  $(-1.96, 1.96)$ , indicando que no se puede rechazar la hipótesis nula al 5% de significancia.

La prueba  $Z$  permite evaluar diferencias entre medias muestrales y poblacionales, especialmente cuando el tamaño de la muestra es grande o la desviación estándar de la población es conocida. Es una prueba útil para decidir si una hipótesis es estadísticamente significativa, aplicando principios de la distribución normal y del Teorema del Límite Central.

#### 4.6 Varianza

La varianza mide la dispersión o variabilidad de los datos con respecto a la media.

En palabras sencillas:

Es como si dijeras: “¿Qué tan lejos están, en promedio, los datos del promedio, pero elevando esas distancias al cuadrado para evitar números negativos?”

La varianza es una medida de qué tan “extendidos” están los datos en un grupo, pero al cuadrar las diferencias, da más peso a los valores extremos (los que están muy lejos del promedio).

Dependiendo de si se analiza toda una población o solo una muestra de esa población, utilizamos dos tipos de varianza:

1. Varianza poblacional ( $\sigma^2$ ): Se utiliza cuando tenemos todos los datos de una población completa. La fórmula es:

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

donde:

- $x_i$  es cada valor del conjunto de datos.
  - $\mu$  es la media de la población.
  - $N$  es el tamaño de la población.
2. Varianza muestral ( $s^2$ ): Se utiliza cuando trabajamos con una muestra de la población. La fórmula para la varianza muestral incluye una corrección, conocida como corrección de Bessel, que divide entre  $n - 1$  en lugar de  $n$  para ajustar el sesgo en muestras pequeñas. La fórmula es:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

donde:

- $x_i$  es cada valor del conjunto de datos muestrales.
- $\bar{x}$  es la media de la muestra.
- $n$  es el tamaño de la muestra.

#### 4.6.1 Utilidad de la varianza

La varianza es útil para:

- Medir la dispersión: Nos permite entender cuán dispersos o cercanos están los datos con respecto a su media.
- Comparar variabilidad: Comparar la varianza entre diferentes conjuntos de datos para determinar cuál tiene más variabilidad.
- Describir la incertidumbre: En estadísticas inferenciales, la varianza es una medida clave para entender la precisión y el comportamiento de los datos.

Ejemplo varianza poblacional y muestral

Supongamos que estamos interesados en la estatura de los estudiantes de una escuela. Si tenemos los datos completos de toda la población de estudiantes, calculamos la varianza poblacional. Si tomamos una muestra de estudiantes, usamos la varianza muestral.

Datos

Para ilustrar, consideremos los siguientes datos de estatura (en cm) de 5 estudiantes:

160, 165, 170, 175, 180

#### 4.6.2 Varianza Poblacional

Dado que tenemos la estatura de todos los estudiantes, calculamos la varianza poblacional.

1. Calcular la media poblacional ( $\mu$ ):

$$\mu = \frac{160 + 165 + 170 + 175 + 180}{5} = 170$$

2. Calcular las desviaciones al cuadrado:

- $(160 - 170)^2 = 100$
- $(165 - 170)^2 = 25$
- $(170 - 170)^2 = 0$
- $(175 - 170)^2 = 25$
- $(180 - 170)^2 = 100$

3. Calcular la varianza poblacional ( $\sigma^2$ ):

$$\sigma^2 = \frac{100 + 25 + 0 + 25 + 100}{5} = \frac{250}{5} = 50$$

La varianza poblacional es 50 cm<sup>2</sup>.

### 4.6.3 Varianza Muestral

Ahora, supongamos que solo tomamos una muestra de 4 estudiantes: 160, 165, 170, 175.

1. Calcular la media muestral ( $\bar{x}$ ):

$$\bar{x} = \frac{160 + 165 + 170 + 175}{4} = 167,5$$

2. Calcular las desviaciones al cuadrado:

- $(160 - 167,5)^2 = 56,25$
- $(165 - 167,5)^2 = 6,25$
- $(170 - 167,5)^2 = 6,25$
- $(175 - 167,5)^2 = 56,25$

3. Calcular la varianza muestral ( $s^2$ ):

$$s^2 = \frac{56,25 + 6,25 + 6,25 + 56,25}{4 - 1} = \frac{125}{3} = 41,67$$

La varianza muestral es 41.67 cm<sup>2</sup>.

#### Conclusión

- La varianza poblacional nos permite conocer la dispersión de todos los datos en la población.
- La varianza muestral corrige el sesgo de subestimar la varianza cuando usamos una muestra.

### 4.7 Desviación Estándar

La desviación estándar es una medida que nos dice qué tan dispersos o diferentes están los datos en un grupo con respecto a su promedio (media).

En palabras sencillas: Es como medir qué tanto se “alejan” los valores individuales del promedio.

#### Interpretación básica

- Si la desviación estándar es pequeña, significa que la mayoría de los datos están cerca del promedio.
- Si la desviación estándar es grande, significa que los datos están más separados y hay más variación.

#### Ejemplo práctico

Supongamos que medimos el peso de 5 cajas de cereal:

500g, 505g, 495g, 502g, 498g.

Aquí, los pesos están bastante cerca del promedio (500g), así que la desviación estándar será baja.

Ahora, si los pesos fueran:

450g, 550g, 400g, 600g, 500g,

hay mucha más diferencia entre los valores y el promedio, por lo que la desviación estándar será alta.

En resumen, la desviación estándar nos ayuda a entender si los datos son consistentes (poco dispersos) o muy variados (muy dispersos).

La desviación estándar tiene dos variantes dependiendo de si estamos analizando una población completa o una muestra:

1. Desviación estándar poblacional ( $\sigma$ ): Se utiliza cuando tenemos todos los datos de una población completa. Es la raíz cuadrada de la varianza poblacional:

$$\sigma = \sqrt{\frac{\sum_{i=1}^N (x_i - \mu)^2}{N}}$$

donde:

- $x_i$  es cada valor del conjunto de datos.
  - $\mu$  es la media de la población.
  - $N$  es el tamaño de la población.
2. Desviación estándar muestral ( $s$ ): Se utiliza cuando trabajamos con una muestra de la población. Es la raíz cuadrada de la varianza muestral, con la corrección por  $n - 1$ :

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

donde:

- $x_i$  es cada valor del conjunto de datos muestrales.
- $\bar{x}$  es la media de la muestra.
- $n$  es el tamaño de la muestra.

#### 4.7.1 Utilidad de la desviación estándar

La desviación estándar es útil para:

- Medir la dispersión: Nos indica cuánto varían los datos respecto a su media, lo que nos permite entender la consistencia o variabilidad de los datos.
- Comparar conjuntos de datos: Con la desviación estándar podemos comparar la variabilidad de diferentes conjuntos de datos, independientemente de sus medias.
- Interpretar la distribución de los datos: En muchas distribuciones (como la normal), una mayor parte de los datos se encuentra dentro de una, dos o tres desviaciones estándar de la media.

Ejemplo: Desviación Estándar Poblacional y Muestral

Supongamos que estamos interesados en la cantidad de horas que un grupo de estudiantes dedica al estudio en una semana. Si tenemos los datos de toda la población, calculamos la desviación estándar poblacional. Si solo contamos con una muestra de estudiantes, usamos la desviación estándar muestral.

Datos

Consideremos las horas de estudio (en horas) de 5 estudiantes:

10, 12, 14, 16, 18

#### 4.7.2 Desviación Estándar Poblacional

Dado que tenemos todos los datos, calculamos la desviación estándar poblacional.

1. Calcular la media poblacional ( $\mu$ ):

$$\mu = \frac{10 + 12 + 14 + 16 + 18}{5} = 14$$

2. Calcular las desviaciones al cuadrado:

- $(10 - 14)^2 = 16$
- $(12 - 14)^2 = 4$
- $(14 - 14)^2 = 0$
- $(16 - 14)^2 = 4$
- $(18 - 14)^2 = 16$

3. Calcular la varianza poblacional:

$$\sigma^2 = \frac{16 + 4 + 0 + 4 + 16}{5} = \frac{40}{5} = 8$$

4. Calcular la desviación estándar poblacional:

$$\sigma = \sqrt{8} = 2,83$$

La desviación estándar poblacional es 2.83 horas.

#### 4.7.3 Desviación Estándar Muestral

Ahora, supongamos que solo tenemos una muestra de 4 estudiantes: 10, 12, 14, 16.

1. Calcular la media muestral ( $\bar{x}$ ):

$$\bar{x} = \frac{10 + 12 + 14 + 16}{4} = 13$$

2. Calcular las desviaciones al cuadrado:

- $(10 - 13)^2 = 9$
- $(12 - 13)^2 = 1$
- $(14 - 13)^2 = 1$
- $(16 - 13)^2 = 9$

3. Calcular la varianza muestral:

$$s^2 = \frac{9 + 1 + 1 + 9}{4 - 1} = \frac{20}{3} = 6,67$$

4. Calcular la desviación estándar muestral:

$$s = \sqrt{6,67} = 2,58$$

La desviación estándar muestral es 2.58 horas.

## Conclusión

- La desviación estándar poblacional mide la dispersión de todos los datos en una población.
- La desviación estándar muestral ajusta la varianza de la muestra para corregir el sesgo en muestras pequeñas, proporcionando una mejor estimación de la variabilidad de la población.

## 4.8 Ejemplo detallado de la Prueba Z para la Media Poblacional

## Prueba de Hipótesis para el Peso de Empaques de Cereal

Una empresa asegura que el peso promedio de un empaque de cereal es de  $\mu = 500$  gramos, con una desviación estándar conocida de  $\sigma = 15$  gramos.

Objetivo: Determinar si el peso promedio de los empaques es significativamente diferente al valor declarado de 500 gramos, usando una prueba Z para la media con un nivel de significancia de  $\alpha = 0,05$ .

## Planteamiento de las Hipótesis

- Hipótesis nula ( $H_0$ ): El peso promedio es igual a 500 gramos.  $H_0 : \mu = 500$
- Hipótesis alternativa ( $H_1$ ): El peso promedio es diferente de 500 gramos.  $H_1 : \mu \neq 500$

Esta es una prueba bilateral porque evaluamos diferencias en ambas direcciones (mayor o menor a 500 gramos).

## Datos de la Muestra

| Peso (g) |
|----------|----------|----------|----------|----------|----------|
| 490      | 495      | 500      | 487      | 498      | 492      |
| 485      | 496      | 502      | 493      | 490      | 488      |
| 491      | 494      | 497      | 486      | 501      | 490      |
| 489      | 493      | 496      | 487      | 495      | 492      |
| 488      | 491      | 498      | 485      | 494      | 489      |
| 487      | 500      | 492      | 490      | 496      | 488      |

Cuadro 4.3: Pesos en gramos de los empaques de cereal.

## Cálculo de la Media Muestral

$$\sum X_i = 490 + 495 + 500 + 487 + \dots + 488 = 17784$$

$$\bar{X} = \frac{\sum X_i}{n} = \frac{17784}{36} = 494 \text{ g.}$$

## Cálculo del Estadístico Z

La fórmula para calcular el valor Z en una prueba de hipótesis es:

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

donde:

- $\bar{X}$ : media muestral,
- $\mu$ : media poblacional bajo la hipótesis nula,
- $\sigma$ : desviación estándar de la población,
- $n$ : tamaño de la muestra.

Calculamos Z:

$$Z = \frac{\bar{X} - \mu}{\text{Error estándar}} = \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{494 - 500}{2,5} = \frac{-6}{2,5} = -2,4$$

## Regla de Decisión

- Valores críticos para  $\alpha = 0,05$  en una prueba bilateral:  $Z_{\text{crítico}} = \pm 1,96$ .
- Como  $Z = -2,4$ , que está fuera del rango  $-1,96 \leq Z \leq 1,96$ , rechazamos  $H_0$ .

$\alpha$	Valor Crítico Z (Dos Colas)
0.10	$\pm 1,645$
0.05	$\pm 1,960$
0.01	$\pm 2,576$
0.001	$\pm 3,291$

Cuadro 4.4: Valores Críticos Z para la Prueba de Dos Colas

## Conclusión

Hay suficiente evidencia estadística para afirmar que el peso promedio de los empaques de cereal difiere significativamente del valor declarado de 500 gramos.

## Resumen de Resultados

---

Parámetro	Valor
Media poblacional ( $\mu$ )	500 g
Media muestral ( $\bar{X}$ )	494 g
Desviación estándar ( $s$ )	15 g
Tamaño de la muestra ( $n$ )	36
Error estándar	2.5 g
Estadístico Z	-2.4
Valores críticos Z	$\pm 1,96$
Decisión	Rechazar $H_0$

---

Cuadro 4.5: Resumen de los cálculos y resultados.





## Bibliografía

- [1] Devore, J. L. (2012). Probabilidad y estadística para ingeniería y ciencias. Internacional Thompson.
- [2] Hildebrand, D. K., & Ott, L. R. (1997). Estadística aplicada a la administración y la economía. Addison-Wesley Iberoamericana.
- [3] Canavos, G. (1988). Probabilidad y estadística. Nueva York: McGraw-Hill.
- [4] Spiegel, M. R. (1999). Estadística. México: McGraw-Hill.