

# **Informe Técnico Completo**

Pipeline Bioestadístico para la Clasificación de  
Sedentarismo mediante Lógica Difusa y Clustering

Perspectiva Bioestadística, Clínica y Computacional

Luis Ángel Martínez

Universidad Autónoma de Chihuahua  
Facultad de Medicina y Ciencias Biomédicas

Programa de Maestría en Ciencias de la Salud

22 de octubre de 2025

## Resumen

El presente informe técnico documenta de manera exhaustiva el pipeline bioestadístico desarrollado para la clasificación objetiva del sedentarismo semanal utilizando datos biométricos de dispositivos wearables (Apple Watch). Este proyecto representa un estudio longitudinal con  $N = 10$  participantes (5M/5H) que generaron 1,337 semanas válidas de datos continuos.

El pipeline integra tres perspectivas complementarias: **bioestadística** (modelado probabilístico robusto, reducción dimensional, clustering, validación), **clínica** (normalización antropométrica, interpretación fisiológica de variables derivadas, relevancia para ciencias del ejercicio), y **computacional** (arquitectura modular en Python, estrategias de imputación jerárquica, optimización de hiperparámetros).

Metodológicamente, el estudio pivotó de un enfoque supervisado inicial (predicción de Calidad de Vida mediante Redes Neuronales Artificiales, invalidado empíricamente) a un paradigma *data-driven* dual: (1) descubrimiento de patrones mediante clustering no supervisado (K-Means,  $K = 2$ , Silhouette= 0,232), empleado como **Verdad Operativa (GO)**, y (2) construcción de un Sistema de Inferencia Difusa Mamdani interpretable con 5 reglas expertas, validado contra la GO con  $F1 = 0,840$ , Recall= 0,976, MCC= 0,294.

Cada fase del pipeline se presenta bajo el marco riguroso de los **6 pasos del análisis estadístico**: planteamiento de hipótesis, selección del estadístico, regla de decisión, cálculos, decisión estadística y conclusión. Se incluyen ecuaciones matemáticas formales, pseudocódigo, referencias a figuras y tablas, y una justificación detallada de la decisión metodológica de *no* emplear un split Train/Test 80/20, reemplazado por validación cruzada Leave-One-User-Out (LOUO) y análisis de sensibilidad.

**Palabras clave:** Sedentarismo, Wearables, Apple Watch, Lógica Difusa, Clustering, K-Means, Imputación Jerárquica, Ingeniería de Características, Validación Cruzada, Python.

# Índice general

<b>1. Planteamiento del Problema e Hipótesis Inicial</b>	<b>5</b>
1.1. Contexto Epidemiológico y Clínico . . . . .	5
1.2. Hipótesis Inicial y Objetivo Primario . . . . .	5
1.2.1. Objetivo Primario (Fase Inicial) . . . . .	6
1.3. Marco de los 6 Pasos: Planteamiento . . . . .	6
<b>2. Selección del Dispositivo Wearable y Diseño de la Cohorte</b>	<b>7</b>
2.1. Evaluación de Dispositivos Wearables . . . . .	7
2.1.1. Criterios de Selección . . . . .	7
2.1.2. Análisis Comparativo . . . . .	7
2.2. Diseño de la Cohorte . . . . .	8
2.2.1. Tamaño Muestral y Justificación . . . . .	8
2.2.2. Criterios de Inclusión/Exclusión . . . . .	9
<b>3. Protocolo de Convocatoria, Recepción y Preprocesamiento de Datos</b>	<b>10</b>
3.1. Protocolo de Recolección de Datos . . . . .	10
3.1.1. Diseño del Protocolo . . . . .	10
3.1.2. Estructura de Datos Crudos . . . . .	10
3.2. Pipeline de Preprocesamiento . . . . .	11
3.2.1. Conversión XML → CSV . . . . .	11
3.2.2. Auditoría de Calidad de Datos . . . . .	12
<b>4. Análisis Exploratorio de Datos (EDA) y Validación del SF-36</b>	<b>13</b>
4.1. Caracterización de Variables Biométricas . . . . .	13
4.1.1. Tipología y Distribuciones . . . . .	13
4.1.2. Gráficos Exploratorios . . . . .	14
4.2. Validación Psicométrica del SF-36 . . . . .	14
4.2.1. Estructura del Cuestionario . . . . .	14
<b>5. Pivote Metodológico: Del Enfoque Supervisado al Data-Driven</b>	<b>16</b>
5.1. Análisis de Correlación SF-36 vs Biométricos . . . . .	16
5.1.1. Hipótesis y Pruebas Iniciales . . . . .	16
5.2. Modelado con Redes Neuronales Artificiales (ANN) . . . . .	17
5.2.1. Arquitectura y Entrenamiento . . . . .	17
5.3. Reformulación: Nuevo Enfoque Data-Driven . . . . .	18

5.3.1. Nueva Hipótesis . . . . .	18
<b>6. Estrategia de Imputación Jerárquica para Datos Faltantes</b>	<b>20</b>
6.1. Diagnóstico de Missingness . . . . .	20
6.1.1. Mecanismos de Datos Faltantes . . . . .	20
6.2. Estrategia de Imputación Jerárquica . . . . .	21
6.2.1. Principios de Diseño . . . . .	21
6.2.2. Algoritmo de Imputación . . . . .	22
6.2.3. Resultados de Imputación . . . . .	23
<b>7. Ingeniería de Características: Variables Derivadas con Normalización Antropométrica</b>	<b>24</b>
7.1. Problema de Comparabilidad Inter-Sujeto . . . . .	24
7.1.1. Heterogeneidad Antropométrica . . . . .	24
7.2. Variable 1: Actividad Relativa . . . . .	25
7.2.1. Definición y Justificación . . . . .	25
7.2.2. Distribución y Validación . . . . .	25
7.3. Variable 2: Superávit Calórico Basal . . . . .	26
7.3.1. Cálculo de TMB . . . . .	26
7.3.2. Definición de Superávit . . . . .	26
7.4. Variables 3 y 4: Perfiles Cardiovasculares . . . . .	26
7.4.1. Delta Cardíaco . . . . .	26
7.4.2. HRV SDNN . . . . .	27
<b>8. Agregación Temporal y Análisis Dual de Variabilidad</b>	<b>28</b>
8.1. Justificación de la Agregación Semanal . . . . .	28
8.1.1. Ventana de Agregación . . . . .	28
8.2. Estadísticos Calculados por Semana . . . . .	28
8.3. Análisis Dual de Variabilidad . . . . .	29
8.3.1. Definición de Variabilidad Observada vs Operativa . . . . .	29
8.3.2. Comparación Observada vs Operativa . . . . .	29
8.3.3. Gráficos de Variabilidad . . . . .	30
8.4. Agregación Semanal: Resultados Finales . . . . .	30
<b>9. Análisis de Correlación, Multicolinealidad y Reducción Dimensional (PCA)</b>	<b>32</b>
9.1. Análisis de Correlación entre Variables Semanales . . . . .	32
9.1.1. Matriz de Correlación . . . . .	32
9.2. Análisis de Multicolinealidad (VIF) . . . . .	33
9.2.1. Factor de Inflación de la Varianza . . . . .	33
9.3. Análisis de Componentes Principales (PCA) . . . . .	35
9.3.1. Reducción Dimensional y Visualización . . . . .	35
<b>10. Clustering No Supervisado: Verdad Operativa (K-Means, K=2)</b>	<b>38</b>
10.1. Justificación del Clustering como Verdad Operativa . . . . .	38
10.1.1. Selección del Algoritmo . . . . .	38
10.2. Barrido de $K$ (K-Sweep) y Selección del Número Óptimo de Clusters . . . . .	39

10.3. Perfiles de Cluster: Análisis Estadístico Detallado . . . . .	40
10.3.1. Asignación de Etiquetas Clínicas . . . . .	40
10.3.2. Estadísticos Descriptivos por Cluster . . . . .	40
10.3.3. Pruebas de Comparación Estadística . . . . .	41
<b>11. Sistema de Inferencia Difusa Mamdani</b>	<b>43</b>
11.1. Diseño del Sistema de Inferencia Difusa . . . . .	43
11.1.1. Arquitectura General . . . . .	43
11.2. Funciones de Pertenencia (Membership Functions) . . . . .	44
11.2.1. Diseño de MF Triangulares Basadas en Percentiles . . . . .	44
11.3. Base de Reglas Difusas . . . . .	46
11.3.1. Reglas Clínicas IF-THEN . . . . .	46
11.3.2. Formalización Matricial . . . . .	47
11.4. Proceso de Inferencia Mamdani . . . . .	47
11.4.1. Paso 1: Fuzzificación . . . . .	47
11.4.2. Paso 2: Activación de Reglas (AND = mínimo) . . . . .	47
11.4.3. Paso 3: Agregación . . . . .	48
11.4.4. Paso 4: Defuzzificación (Centroide Discreto) . . . . .	48
11.4.5. Paso 5: Binarización . . . . .	48
<b>12. Validación Cruzada y Análisis de Robustez</b>	<b>49</b>
12.1. Validación por Concordancia: Fuzzy vs Clustering . . . . .	49
12.1.1. Métricas de Desempeño . . . . .	49
12.2. Validación Cruzada Leave-One-User-Out (LOUO) . . . . .	51
12.2.1. Justificación de LOUO . . . . .	51
12.3. Análisis de Sensibilidad . . . . .	52
12.3.1. Sensibilidad al Umbral $\tau$ . . . . .	52
12.3.2. Sensibilidad a Parámetros de MF . . . . .	52
12.4. Análisis de Robustez: Modelo 4V vs Modelo 2V . . . . .	53
12.4.1. Motivación del Análisis . . . . .	53
<b>13. Justificación Metodológica: Por Qué NO Split Train/Test 80/20</b>	<b>55</b>
13.1. Problemática del Split Tradicional en Datos Longitudinales . . . . .	55
13.2. Razón 1: Fuga Temporal (Temporal Leakage) . . . . .	56
13.2.1. Naturaleza de los Datos . . . . .	56
13.3. Razón 2: Insuficiencia de Poder Estadístico . . . . .	57
13.3.1. Split por Usuario vs Split por Semanas . . . . .	57
13.4. Razón 3: Objetivo Descriptivo vs Predictivo . . . . .	58
13.4.1. Naturaleza del Estudio . . . . .	58
13.5. Alternativas Metodológicas Implementadas . . . . .	58
13.5.1. Estrategia de Validación Adoptada . . . . .	58
13.5.2. Leave-One-User-Out (LOUO) Cross-Validation . . . . .	59
13.6. Resumen de Defensa Metodológica . . . . .	59

# Capítulo 1

## Planteamiento del Problema e Hipótesis Inicial

### 1.1. Contexto Epidemiológico y Clínico

El comportamiento sedentario (CS), definido por la Organización Mundial de la Salud como cualquier actividad con gasto energético  $\leq 1,5$  METs en posición sentada o reclinada durante horas de vigilia, constituye un factor de riesgo independiente para enfermedades crónicas no transmisibles (ECNT), incluyendo obesidad, diabetes tipo 2, enfermedad cardiovascular y ciertos tipos de cáncer [1].

La medición objetiva del CS mediante acelerometría triaxial en dispositivos wearables de consumo masivo (e.g., Apple Watch, Fitbit, Garmin) ha revolucionado la epidemiología del comportamiento, permitiendo cuantificar patrones de actividad física en condiciones de “vida libre” con alta resolución temporal ( $\geq 1$  Hz) y sin el sesgo de auto-reporte característico de cuestionarios.

### 1.2. Hipótesis Inicial y Objetivo Primario

#### Paso 1: Planteamiento de Hipótesis

##### **Hipótesis $H_0$ (inicial, posteriormente rechazada):**

Existe una relación inversa, lineal y medible entre el comportamiento sedentario objetivo (CS\_obj), cuantificado mediante métricas derivadas de acelerometría y fotopleismografía (PPG) del Apple Watch, y la percepción subjetiva de Calidad de Vida Relacionada con la Salud (CVRS), evaluada mediante el cuestionario SF-36.

Formalmente:

$$CVRS_{SF36} = f(CS_{obj}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \quad (1.1)$$

donde  $f$  sería una función lineal o no lineal modelable mediante Redes Neuronales Artificiales (ANN).

### 1.2.1. Objetivo Primario (Fase Inicial)

Desarrollar un modelo predictivo (ANN) capaz de cuantificar la CVRS a partir de métricas biométricas continuas, con  $R^2 \geq 0,70$  y MAE  $\leq 10$  puntos en escala SF-36.

## 1.3. Marco de los 6 Pasos: Planteamiento

### Paso 2: Selección del Estadístico/Método

#### Selección del método:

Se propuso inicialmente un análisis correlacional (Pearson/Spearman) seguido de modelado supervisado mediante ANN (arquitectura feedforward, activación ReLU, optimizador Adam).

### Paso 3: Regla de Decisión

#### Regla de decisión:

Si  $|r| \geq 0,60$  (correlación fuerte) y el modelo ANN alcanza  $R^2 \geq 0,70$  en validación cruzada 5-fold, se aceptará la hipótesis de relación cuantificable.

### Paso 5: Decisión Estadística

#### Decisión preliminar:

Se decidió proceder con un diseño longitudinal que recolectaría datos biométricos continuos (Apple Watch) y evaluaciones periódicas del SF-36 para probar esta correlación.

### Paso 6: Conclusión

#### Conclusión del planteamiento:

Existía suficiente justificación teórica (revisión de literatura: correlaciones reportadas entre actividad física y CVRS en el rango  $r = 0,30 - 0,50$ ) para explorar esta vía, aunque con la precaución de que la relación podría ser más compleja de lo anticipado.

## Capítulo 2

# Selección del Dispositivo Wearable y Diseño de la Cohorte

### 2.1. Evaluación de Dispositivos Wearables

#### 2.1.1. Criterios de Selección

##### Paso 1: Planteamiento de Hipótesis

###### Problema/Hipótesis:

Necesitábamos un dispositivo wearable que cumpliera simultáneamente:

- Alta penetración de mercado (facilitar reclutamiento BYOD)
- Sensores validados: acelerómetro 3-ejes ( $\geq 50$  Hz), PPG para FC/VFC
- Plataforma de exportación de datos crudos o agregados
- Consistencia inter-versión (minimizar heterogeneidad instrumental)

Hipótesis: El Apple Watch, por su ecosistema cerrado y validaciones previas en literatura (Stahl et al., 2016; Shcherbina et al., 2017), sería la opción preferente.

#### 2.1.2. Análisis Comparativo

Tabla 2.1: Matriz de Decisión: Comparación de Dispositivos Wearables

Criterio	Apple Watch	Fitbit	Garmin	Mi Band
Penetración México	Alta	Media	Media-Baja	Alta
Sensores validados	Sí	Sí	Sí	Parcial
Exportación datos	HealthKit (XML)	API limitada	Garmin Connect	Propietaria
Consistencia HW	Alta	Media	Alta	Baja
Costo promedio (USD)	300-800	100-300	250-700	30-50
Score ponderado	9.2	7.5	7.8	5.1



### Paso 2: Selección del Estadístico/Método

#### Método de evaluación:

Matriz de decisión multicriterio con pesos asignados según importancia para el estudio:

- Validez de sensores: 35 %
- Exportabilidad de datos: 30 %
- Consistencia: 20 %
- Penetración: 15 %

### Paso 5: Decisión Estadística

#### Decisión:

Se seleccionó el **Apple Watch** (Series 3 o superior) como dispositivo estándar del estudio, adoptando un enfoque *Bring Your Own Device* (BYOD) para maximizar adherencia y minimizar el efecto Hawthorne.

## 2.2. Diseño de la Cohorte

### 2.2.1. Tamaño Muestral y Justificación

#### Paso 1: Planteamiento de Hipótesis

##### Planteamiento:

Dada la naturaleza longitudinal del estudio (objetivo: capturar variabilidad intra-sujeto durante  $\geq 12$  semanas), el tamaño muestral  $N$  se justificó por:

$$n_{\text{observaciones}} = N_{\text{sujetos}} \times T_{\text{semanas}} \geq 1000 \quad (2.1)$$

Con  $N = 10$  y  $T \approx 130$  semanas (promedio), se alcanzarían  $\approx 1300$  observaciones semanales, suficiente para:

- Modelado de clustering con  $n/K \geq 500$  por grupo ( $K = 2$ )
- Optimización de hiperparámetros del sistema difuso
- Validación cruzada Leave-One-Subject-Out

2.2.2. Criterios de Inclusión/Exclusión

Tabla 2.2: Criterios de Elegibilidad de Participantes

Criterio	Inclu
Edad	18-65
Dispositivo	Propietario Apple Watch Series
Uso previo	$\geq 6$ meses cont
Estado de salud	Ambulatorio, sin limitac
Consentimiento	Informado por es
Datos exportables	$\geq 80\%$ días con c

Paso 4: Cálculos

**Cálculos de factibilidad:**

Se convocó a 15 candidatos, de los cuales:

- 12 cumplieron criterios de inclusión
- 10 completaron el protocolo (2 abandonos por causas no relacionadas)
- Distribución final: 5 hombres, 5 mujeres
- Edad:  $\bar{x} = 32,4$  años,  $s = 8,7$  años
- IMC:  $\bar{x} = 26,1$  kg/m<sup>2</sup>,  $s = 4,2$  kg/m<sup>2</sup>

Paso 6: Conclusión

**Conclusión metodológica:**

Aunque no representativa poblacionalmente (muestra de conveniencia), la cohorte de  $N = 10$  permite un análisis longitudinal profundo con potencia estadística adecuada para el descubrimiento de patrones intra-sujeto y validación de sistemas expertos interpretativos (objetivo secundario tras el pivote metodológico).

## Capítulo 3

# Protocolo de Convocatoria, Recepción y Preprocesamiento de Datos

### 3.1. Protocolo de Recolección de Datos

#### 3.1.1. Diseño del Protocolo

##### Paso 1: Planteamiento de Hipótesis

###### Planteamiento:

Para garantizar la integridad, trazabilidad y ética de los datos biométricos sensibles, se diseñó un protocolo estandarizado que incluye:

1. Consentimiento informado (aprobación comité ética institucional)
2. Instrucciones de exportación (HealthKit → archivo `export.zip`)
3. Aplicación del cuestionario SF-36 (versión mexicana validada)
4. Anonimización inmediata (códigos: u1, u2, ..., u10)
5. Almacenamiento seguro (servidor institucional, encriptación AES-256)

#### 3.1.2. Estructura de Datos Crudos

Los datos exportados de Apple Health siguen el esquema XML:

```
1 <HealthData>
2   <Record type="HKQuantityTypeIdentifierStepCount"
3       sourceName="Apple Watch de Luis"
4       value="1245"
5       unit="count"
6       startDate="2023-10-22 08:15:00"
7       endDate="2023-10-22 08:16:00"/>
8   ...
9 </HealthData>
```

Listing 3.1: Estructura XML de Apple Health Export

## 3.2. Pipeline de Preprocesamiento

### 3.2.1. Conversión XML → CSV

#### Paso 2: Selección del Estadístico/Método

##### Método:

Parseo XML mediante `ElementTree` (Python), con transformaciones:

- Filtrado por `sourceName` (solo datos Apple Watch, excluir iPhone)
- Conversión de timestamps a zona horaria local (UTC-6, Chihuahua)
- Agregación a nivel diario (suma/media según métrica)

---

#### Algorithm 1 Preprocesamiento XML a CSV Diario

---

```

1: Input: export.zip por participante
2: Output: DB_u{id}.csv con columnas [fecha, pasos, calorías, fc_reposo, hrv_sdnm, ...]
3:
4: procedure PARSEXML(xml_file, user_id)
5:   tree ← parse(xml_file)
6:   records ← tree.findall("Record")
7:   df ← empty_dataframe()
8:   for record in records do
9:     if record.sourceName contains ".Apple Watch" then
10:       type ← record.type
11:       value ← record.value
12:       date ← record.startDate.date()
13:       df.append([date, type, value])
14:     end if
15:   end for
16:   df_pivot ← df.pivot(index=date, columns=type, values=value)
17:   df_pivot.to_csv(f"DB_u{user_id}.csv")
18: end procedure

```

---

#### Paso 4: Cálculos

##### Cálculos de agregación:

Para cada usuario y día:

$$\text{Pasos}_{\text{día}} = \sum_{t=0}^{23:59} \text{StepCount}(t) \quad (3.1)$$

$$\text{FC}_{\text{reposo}} = \min\{\text{HeartRate}(t) : t \in [02 : 00, 05 : 00]\} \quad (3.2)$$

$$\text{HRV\_SDNN}_{\text{día}} = \text{mean}\{\text{SDNN}(t) : t \in [00 : 00, 23 : 59]\} \quad (3.3)$$

### 3.2.2. Auditoría de Calidad de Datos

Tabla 3.1: Métricas de Completitud por Usuario (Fase Pre-Imputación)

Usuario	Días totales	Días válidos	Completitud (%)	Missing FC (%)	Missing HRV (%)
u1	900	852	94.7	8.2	15.1
u2	850	801	94.2	9.1	17.1
u3	920	884	96.1	5.4	12.1
...	...	...	...	...	...
u10	880	831	94.4	7.8	14.1
<b>Media</b>	<b>885</b>	<b>838</b>	<b>94.7</b>	<b>7.6</b>	<b>14.1</b>

#### Paso 5: Decisión Estadística

##### Decisión:

La completitud general  $> 94\%$  es aceptable para estudios observacionales de vida libre. Las variables cardiovasculares (FC, HRV) presentan mayor tasa de missin-gness (mecanismo: quitarse el reloj durante sueño/carga), requiriendo estrategia de imputación robusta (Capítulo 6).

# Capítulo 4

## Análisis Exploratorio de Datos (EDA) y Validación del SF-36

### 4.1. Caracterización de Variables Biométricas

#### 4.1.1. Tipología y Distribuciones

##### Paso 1: Planteamiento de Hipótesis

###### Hipótesis:

Se esperaba que las variables biométricas diarias presentaran:

- Distribuciones asimétricas (pasos, minutos ejercicio: asimetría positiva)
- Alta variabilidad día-a-día ( $CV > 50\%$ )
- No-normalidad (rechazo de Shapiro-Wilk con  $p < 0,05$ )

##### Paso 2: Selección del Estadístico/Método

###### Métodos aplicados:

- Estadísticos descriptivos robustos: mediana, IQR, MAD
- Pruebas de normalidad: Shapiro-Wilk (si  $n < 5000$ ), Kolmogorov-Smirnov (si  $n \geq 5000$ )
- Visualización: histogramas, Q-Q plots, boxplots por usuario

Tabla 4.1: Estadísticos Descriptivos de Variables Clave (Nivel Diario,  $n = 8,380$  días)

Variable	Media	DE	Mediana	IQR	Min	Max	SW $p$ -valor
Pasos	6,842	4,231	6,120	4,890	0	28,450	$< 0,001$
Calorías activas	385	287	342	298	0	1,892	$< 0,001$
FC reposo (lpm)	58.3	8.7	57.0	10.0	42	92	0,014
HRV SDNN (ms)	52.1	18.4	48.5	22.0	15	128	$< 0,001$
FC caminar (lpm)	95.8	12.3	94.0	15.0	65	145	0,082
Min sedentarios	678	142	702	185	120	1,320	$< 0,001$

### Paso 5: Decisión Estadística

#### Decisión estadística:

Se rechaza la normalidad para todas las variables excepto FC\_caminar ( $p = 0,082$ ).  
Consecuencia: uso obligatorio de métodos no paramétricos o robustos (medianas, bootstrapping, Mann-Whitney U) en análisis posteriores.

## 4.1.2. Gráficos Exploratorios

*Ver Figuras:*

- 4 semestre\_dataset/analisis\_u/histogramas\_variables\_clave.png
- 4 semestre\_dataset/analisis\_u/qqplots\_normalidad.png
- 4 semestre\_dataset/analisis\_u/boxplots\_por\_usuario.png

## 4.2. Validación Psicométrica del SF-36

### 4.2.1. Estructura del Cuestionario

El SF-36 evalúa 8 dimensiones de CVRS mediante 36 ítems:

- Función Física (FF)
- Rol Físico (RF)
- Dolor Corporal (DC)
- Salud General (SG)
- Vitalidad (VT)
- Función Social (FS)
- Rol Emocional (RE)
- Salud Mental (SM)

**Paso 2: Selección del Estadístico/Método****Métrica de fiabilidad:**

Alfa de Cronbach por dimensión, criterio  $\alpha \geq 0,70$  (aceptable).

$$\alpha = \frac{K}{K-1} \left( 1 - \frac{\sum_{i=1}^K \sigma_i^2}{\sigma_{\text{total}}^2} \right) \quad (4.1)$$

donde  $K$  = número de ítems,  $\sigma_i^2$  = varianza del ítem  $i$ .

Tabla 4.2: Fiabilidad del SF-36 en la Cohorte ( $N = 10$ )

Dimensión SF-36	$\alpha$ Cronbach	Varianza	Decisión
Función Física	0.82	145.3	Aceptable
Rol Físico	0.51	0.0	Rechazada (var=0)
Dolor Corporal	0.78	98.7	Aceptable
Salud General	0.73	112.4	Aceptable
Vitalidad	0.64	87.2	Marginal
Función Social	0.71	102.1	Aceptable
Rol Emocional	0.76	118.5	Aceptable
Salud Mental	0.80	134.2	Aceptable

**Paso 5: Decisión Estadística****Decisión crítica:**

La dimensión **Rol Físico** presenta varianza nula (todos los participantes reportaron el mismo valor, efecto techo/suelo), invalidando su uso. Vitalidad ( $\alpha = 0,64$ ) está por debajo del umbral.

**Consecuencia:** Estos problemas psicométricos, sumados a correlaciones débiles con biométricos (siguiente sección), motivaron el rechazo de la hipótesis inicial y el pivote metodológico.

**Paso 6: Conclusión****Conclusión EDA:**

1. Los datos biométricos son ruidosos y no-normales, requiriendo métodos robustos.
2. El SF-36 presenta limitaciones en esta cohorte específica (tamaño, homogeneidad).
3. La alta variabilidad diaria ( $CV > 100\%$  en ejercicio) justifica agregación temporal (semanal) para capturar patrones estables.



## Capítulo 5

# Pivote Metodológico: Del Enfoque Supervisado al Data-Driven

### 5.1. Análisis de Correlación SF-36 vs Biométricos

#### 5.1.1. Hipótesis y Pruebas Iniciales

##### Paso 1: Planteamiento de Hipótesis

###### Hipótesis $H_1$ a probar:

Las métricas biométricas agregadas (media de 4 semanas) correlacionan significativamente ( $|r| \geq 0,60$ ,  $p < 0,01$ ) con los puntajes de CVRS del SF-36.

##### Paso 2: Selección del Estadístico/Método

###### Métodos:

- Correlación de Spearman (datos no-normales)
- Corrección Bonferroni para comparaciones múltiples ( $\alpha^* = 0,05/32 = 0,0016$ )
- Scatter plots con líneas de regresión LOWESS

Tabla 5.1: Matriz de Correlación: Biométricos Agregados vs SF-36 ( $N = 10$ )

	FF	RF	DC	SG	VT	FS	RE	SM
Pasos promedio	0.32	—	0.18	0.41	-0.05	0.27	0.14	0.09
Calorías promedio	0.38	—	0.22	0.45	-0.12	0.31	0.19	0.13
FC reposo promedio	-0.21	—	-0.14	-0.28	0.08	-0.18	-0.11	-0.06
HRV SDNN promedio	0.15	—	0.09	0.24	0.31	0.12	0.08	0.19
Min sedentarios	-0.29	—	-0.16	-0.35	-0.18	-0.24	-0.13	-0.11

*Nota:* RF excluido por varianza nula. Ninguna correlación alcanza  $|r| \geq 0,60$  ni  $p < 0,0016$ .

Paso 5: Decisión Estadística

Decisión estadística:  
Se rechaza  $H_1$ . Las correlaciones observadas son débiles a moderadas ( $0,09 \leq |r| \leq 0,45$ ) y ninguna sobrevive la corrección Bonferroni. La asociación es insuficiente para justificar un modelo predictivo.

## 5.2. Modelado con Redes Neuronales Artificiales (ANN)

### 5.2.1. Arquitectura y Entrenamiento

A pesar de las correlaciones débiles, se procedió a entrenar ANNs como prueba definitiva:

Algorithm 2 Entrenamiento de ANN para CVRS

1: **Input:**  $X \in \mathbb{R}^{10 \times 16}$  (16 features biométricos),  $y \in \mathbb{R}^{10 \times 7}$  (7 dimensiones SF-36 válidas)

2: **Output:** Modelo ANN, métricas de desempeño

3:

4: Arquitectura: [16 inputs]  $\rightarrow$  [32 ReLU]  $\rightarrow$  [16 ReLU]  $\rightarrow$  [7 Linear]

5: Optimizador: Adam ( $\alpha = 0,001$ ,  $\beta_1 = 0,9$ ,  $\beta_2 = 0,999$ )

6: Función de pérdida: MSE

7: Validación cruzada: 5-fold

8: Épocas: 500 con early stopping (patience=50)

Paso 4: Cálculos

Resultados del entrenamiento:

Métrica	Train	Validación	Test	Criterio
$R^2$	0.92	-0.18	-0.34	$\geq 0,70$
MAE	5.2	18.7	21.3	$\leq 10$
RMSE	7.8	24.1	27.9	$\leq 15$

Tabla 5.2: Desempeño del modelo ANN (peor de 20 configuraciones probadas)

Observación crítica:  $R^2$  negativo en validación/test indica que el modelo es *peor que predecir la media*, evidenciando sobreajuste severo y ausencia de relación generalizable.

### Paso 5: Decisión Estadística

#### Decisión metodológica CRÍTICA:

Se rechaza definitivamente la hipótesis inicial y el enfoque supervisado. Las causas identificadas:

1.  $N = 10$  es insuficiente para ANN (regla de oro:  $\geq 10 \times$  parámetros; aquí:  $\approx 1,000$  parámetros)
2. Relación CS-CVRS es multifactorial, confundida por variables psicosociales no capturadas
3. SF-36 carece de sensibilidad a variaciones diarias/semanales de actividad en población joven-adulta sana

## 5.3. Reformulación: Nuevo Enfoque Data-Driven

### 5.3.1. Nueva Hipótesis

#### Paso 1: Planteamiento de Hipótesis

##### Hipótesis $H_2$ (reformulada):

Los datos biométricos contienen patrones latentes que permiten clasificar objetivamente semanas como “alto sedentarismo” vs “bajo sedentarismo”, independientemente de la percepción subjetiva de CVRS.

##### Enfoque dual propuesto:

1. **Descubrimiento empírico:** Clustering no supervisado (K-Means) para identificar grupos naturales en los datos  $\rightarrow$  *Verdad Operativa (GO)*
2. **Sistema experto interpretable:** Lógica Difusa (Mamdani) con reglas basadas en conocimiento fisiológico  $\rightarrow$  *Modelo Clínico*
3. **Validación cruzada:** Concordancia entre ambos métodos independientes

#### Paso 2: Selección del Estadístico/Método

##### Métricas de éxito reformuladas:

- F1-Score  $\geq 0,80$  (balance precisión-recall)
- Matthews Correlation Coefficient (MCC)  $\geq 0,30$  (manejo desbalanceo)
- Interpretabilidad clínica de las reglas difusas

**Paso 6: Conclusión****Conclusión del pivote:**

Este cambio paradigmático transforma el estudio de *predictivo supervisado* a *descriptivo-clasificadorio data-driven*, más apropiado para la naturaleza exploratoria de los datos y el tamaño muestral. Los capítulos siguientes desarrollan este nuevo enfoque.

# Capítulo 6

## Estrategia de Imputación Jerárquica para Datos Faltantes

### 6.1. Diagnóstico de Missingness

#### 6.1.1. Mecanismos de Datos Faltantes

##### Paso 1: Planteamiento de Hipótesis

###### Hipótesis sobre mecanismos:

Los datos faltantes en wearables no son MCAR (Missing Completely At Random), sino:

- **MAR (Missing At Random)**: FC/HRV ausentes durante actividades acuáticas (no resistance device)
- **MNAR (Missing Not At Random)**: Dispositivo quitado intencionalmente durante eventos sedentarios prolongados (e.g., cine, sueño extendido)

##### Paso 2: Selección del Estadístico/Método

###### Pruebas aplicadas:

- Test de Little MCAR:  $\chi^2 = 487,3$ ,  $p < 0,001 \rightarrow$  Rechazo MCAR
- Patrones de missingness visualizados con matrices de co-ocurrencia
- Análisis temporal: ACF/PACF de indicadores de missingness

*Ver Figuras:*

- 4 semestre\_dataset/analisis\_u/missingness\_y\_acf/missingness\_matrix\_u1.png
- 4 semestre\_dataset/analisis\_u/missingness\_y\_acf/acf\_plots/acf\_u1.png
- 4 semestre\_dataset/analisis\_u/missingness\_y\_acf/pacf\_plots/pacf\_u1.png

## 6.2. Estrategia de Imputación Jerárquica

### 6.2.1. Principios de Diseño

1. **Sin fuga temporal:** Imputación *forward-only* (día  $t$  usa solo info  $\leq t - 1$ )
2. **Plausibilidad fisiológica:** Valores imputados dentro de rangos clínicos
3. **Jerarquía de métodos:** De específico a general
4. **Transparencia:** Marcar columnas con sufijo `_imp` y registrar tasa

### 6.2.2. Algoritmo de Imputación

---

**Algorithm 3** Imputación Jerárquica para Variables Cardiovasculares

---

```

1: Input: DataFrame diario con columnas [fecha, FC_caminar, FC_reposo,
   HRV_SDNN, ...]
2: Output: DataFrame con valores imputados y flags
3:
4: for variable in [FC_caminar, FC_reposo, HRV_SDNN] do
5:   for row_idx in missing_indices(variable) do
6:     usuario  $\leftarrow$  row_idx.usuario
7:     fecha  $\leftarrow$  row_idx.fecha
8:
9:     // Método 1: Media móvil 7 días previos
10:    ventana  $\leftarrow$  [fecha-7, fecha-1]
11:    if count(ventana)  $\geq$  4 then
12:      impute median(ventana) ▷ Robusto a outliers
13:      continue
14:    end if
15:
16:    // Método 2: Media del mismo día de semana (último mes)
17:    mismo_dia  $\leftarrow$  filter(fecha.weekday == dia_semana, fecha  $\in$  [fecha-28,
    fecha-1])
18:    if count(mismo_dia)  $\geq$  2 then
19:      impute median(mismo_dia)
20:      continue
21:    end if
22:
23:    // Método 3: Mediana histórica del usuario
24:    historico  $\leftarrow$  filter(usuario == usuario, fecha < fecha)
25:    if count(historico)  $\geq$  10 then
26:      impute median(historico)
27:      continue
28:    end if
29:
30:    // Método 4: Estimación por ecuaciones de Tanaka (FC_reposo)
31:    if variable == FC_reposo and edad disponible then
32:      impute  $220 - \text{edad} \times 0,7$  ▷ FC reposo estimado
33:      continue
34:    end if
35:
36:    // Método 5 (último recurso): Mediana global
37:    impute median_global(variable)
38:  end for
39: end for

```

---

### 6.2.3. Resultados de Imputación

Tabla 6.1: Tasa de Imputación por Variable y Método

Variable	Missing (%)	M1 (%)	M2 (%)	M3 (%)	M4 (%)	M5 (%)
FC_caminar	7.6	68.2	21.3	8.9	0.0	1.6
FC_reposo	4.2	72.1	18.7	6.5	2.1	0.6
HRV_SDNN	14.8	61.5	24.8	10.3	0.0	3.4

#### Paso 4: Cálculos

##### Validación de plausibilidad:

Post-imputación, se verificó que todos los valores cumplan:

$$40 \leq FC_{\text{reposo}} \leq 100 \text{ lpm} \quad (6.1)$$

$$60 \leq FC_{\text{caminar}} \leq 160 \text{ lpm} \quad (6.2)$$

$$15 \leq HRV\_SDNN \leq 150 \text{ ms} \quad (6.3)$$

Violaciones detectadas: 3 outliers extremos (0.04 %), reemplazados por mediana del usuario.

#### Paso 5: Decisión Estadística

##### Decisión:

La estrategia jerárquica logró reducir missingness de 14.8 % (HRV) a 0 %, con > 90 % de valores imputados mediante métodos específicos del usuario (M1-M3), garantizando consistencia individual.

#### Paso 6: Conclusión

##### Conclusión:

La imputación jerárquica sin fuga temporal preserva la integridad de series temporales para análisis posteriores (ACF/PACF, agregación semanal). El análisis de variabilidad dual (Capítulo 8) confirmará que la imputación no distorsiona las distribuciones originales.



# Capítulo 7

## Ingeniería de Características: Variables Derivadas con Normalización Antropométrica

### 7.1. Problema de Comparabilidad Inter-Sujeto

#### 7.1.1. Heterogeneidad Antropométrica

##### Paso 1: Planteamiento de Hipótesis

**Problema:**

Variables brutas (pasos, calorías, FC) no son directamente comparables entre individuos con diferente:

- Masa corporal (IMC:  $19.8 - 32.4 \text{ kg/m}^2$  en la cohorte)
- Tasa Metabólica Basal (TMB: función de sexo, edad, peso, altura)
- Tiempo de uso del dispositivo ( $6.2 - 23.8 \text{ h/día}$ )

**Consecuencia:** Un usuario pesado quemará más calorías en reposo que uno liviano; ignorar esto induce sesgo en clustering.

## 7.2. Variable 1: Actividad Relativa

### 7.2.1. Definición y Justificación

#### Paso 2: Selección del Estadístico/Método

##### Derivación matemática:

$$\text{Actividad\_relativa}_{\text{día}} = \frac{\text{Pasos}}{\text{Horas\_con\_datos}} \times \frac{1}{1000} \quad (7.1)$$

Unidades: *kilopasos por hora de monitoreo*

**Justificación clínica:** Normaliza por exposición al dispositivo. Un usuario con 10,000 pasos en 10 horas (1.0 kph) es *más activo* que uno con 10,000 pasos en 20 horas (0.5 kph).

### 7.2.2. Distribución y Validación

Tabla 7.1: Comparación: Pasos Brutos vs Actividad Relativa

Variable	Usuario	Media	DE	CV (%)	Mediana	IQR
Pasos	u1 (IMC 22.1)	8,542	3,921	45.9	8,120	4,650
	u5 (IMC 29.8)	5,234	2,814	53.8	5,010	3,210
	u9 (IMC 24.5)	7,892	3,654	46.3	7,650	4,120
Act_rel (kph)	u1	0.62	0.28	45.2	0.59	0.31
	u5	0.58	0.31	53.4	0.55	0.35
	u9	0.65	0.30	46.2	0.63	0.34

#### Paso 5: Decisión Estadística

##### Decisión:

Actividad\_rel reduce la varianza inter-sujeto atribuible a diferencias en tiempo de uso (CV similar, pero medianas más homogéneas), permitiendo clustering más justo.

## 7.3. Variable 2: Superávit Calórico Basal

### 7.3.1. Cálculo de TMB

#### Paso 2: Selección del Estadístico/Método

##### Ecuación de Harris-Benedict (revisada):

Para hombres:

$$\text{TMB}_h = 88,362 + (13,397 \times \text{peso\_kg}) + (4,799 \times \text{altura\_cm}) - (5,677 \times \text{edad}) \quad (7.2)$$

Para mujeres:

$$\text{TMB}_m = 447,593 + (9,247 \times \text{peso\_kg}) + (3,098 \times \text{altura\_cm}) - (4,330 \times \text{edad}) \quad (7.3)$$

### 7.3.2. Definición de Superávit

$$\text{Superávit\_calórico\_basal}_{\text{día}} = \frac{\text{Calorías\_activas}}{\text{TMB}} \times 100 \% \quad (7.4)$$

#### Interpretación clínica:

- < 20 %: Gasto activo muy bajo (sedentarismo)
- 20 – 50 %: Actividad ligera-moderada
- > 50 %: Actividad vigorosa o deportiva

Tabla 7.2: TMB y Superávit Calórico por Usuario

Usuario	Sexo	IMC	TMB (kcal/día)	Sup. p50 (%)
u1	M	22.1	1,742	28.3
u2	F	24.3	1,521	31.7
u3	M	26.8	1,865	25.9
...	...	...	...	...
u10	F	23.5	1,498	34.2

## 7.4. Variables 3 y 4: Perfiles Cardiovasculares

### 7.4.1. Delta Cardíaco

$$\text{Delta\_cardíaco}_{\text{día}} = \text{FC\_caminar} - \text{FC\_reposo} \quad (7.5)$$

**Relevancia fisiológica:** Mayor delta indica mejor reserva cardiovascular (respuesta rápida del sistema nervioso autónomo a demanda metabólica).

### 7.4.2. HRV SDNN

La Variabilidad de la Frecuencia Cardíaca (HRV), específicamente SDNN (Standard Deviation of NN intervals), es un biomarcador del tono vagal:

- SDNN > 50 ms: Buena modulación autonómica
- SDNN < 30 ms: Posible fatiga, sobreentrenamiento, o estrés crónico

#### Paso 4: Cálculos

**Correlación entre variables derivadas:**

	Act_rel	Sup_cal	HRV	Delta_card
Act_rel	1.00	0.68	0.12	0.24
Sup_cal	0.68	1.00	0.09	0.31
HRV	0.12	0.09	1.00	0.18
Delta_card	0.24	0.31	0.18	1.00

Tabla 7.3: Matriz de Correlación (Spearman,  $n = 8,380$  días)

**Observación:** Correlación moderada Act\_rel – Sup\_cal (esperada: ambas reflejan volumen de actividad), pero baja con variables cardiovasculares, confirmando que capturan dominios distintos.

#### Paso 6: Conclusión

##### Conclusión:

Las 4 variables derivadas son:

1. Antropométricamente normalizadas (comparabilidad)
2. Fisiológicamente interpretables (relevancia clínica)
3. Relativamente independientes ( $r < 0,70$ , evitando multicolinealidad severa)

Estas formarán la base para la agregación semanal (siguiente capítulo) y posterior modelado.

## Capítulo 8

# Agregación Temporal y Análisis Dual de Variabilidad

### 8.1. Justificación de la Agregación Semanal

#### Paso 1: Planteamiento de Hipótesis

##### Hipótesis:

Los datos diarios presentan una variabilidad excesiva ( $CV > 50\%$ ) atribuible a:

- Comportamientos esporádicos (ejercicio intenso 1 día, sedentarismo el siguiente)
- Ruido de medición (errores de sensor, eventos atípicos)
- Ciclos semanales (diferencias fin de semana vs días laborales)

La agregación a nivel semanal (7 días continuos) utilizando estadísticos robustos (mediana, IQR) capturará el *patrón habitual* de comportamiento, reduciendo ruido y mejorando estabilidad para clustering/modelado.

#### 8.1.1. Ventana de Agregación

$$\text{Semana } k : \quad \text{fecha\_inicio} = \text{Lunes}, \quad \text{fecha\_fin} = \text{Domingo} \quad (8.1)$$

**Criterio de validez:** Semana incluida si  $\geq 5$  días tienen datos completos (71 % completitud).

### 8.2. Estadísticos Calculados por Semana

Para cada una de las 4 variables derivadas:

$$x_{p50}^{(k)} = \text{median}\{x_{\text{día}_1}, x_{\text{día}_2}, \dots, x_{\text{día}_7}\} \quad (8.2)$$

$$x_{\text{IQR}}^{(k)} = Q_3(x) - Q_1(x) \quad (8.3)$$

$$x_{p10}^{(k)} = \text{percentil}_{10}(x) \quad (8.4)$$

$$x_{p90}^{(k)} = \text{percentil}_{90}(x) \quad (8.5)$$

Resultado: Dataset semanal con  $n_{\text{semanas}} = 1,337$  (válidas) y 16 features (4 variables  $\times$  4 estadísticos).

## 8.3. Análisis Dual de Variabilidad

### 8.3.1. Definición de Variabilidad Observada vs Operativa

#### Paso 2: Selección del Estadístico/Método

##### Variabilidad Observada (datos crudos, sin imputar):

Cuantifica la fluctuación natural día-a-día medida directamente por el sensor.

$$CV_{\text{obs}}^{(u,v)} = \frac{\sigma_{\text{obs}}(v, u)}{\mu_{\text{obs}}(v, u)} \times 100 \% \quad (8.6)$$

donde  $v$  = variable,  $u$  = usuario.

##### Variabilidad Operativa (datos post-imputación):

Refleja la variabilidad utilizada en el análisis final.

$$CV_{\text{op}}^{(u,v)} = \frac{\sigma_{\text{op}}(v, u)}{\mu_{\text{op}}(v, u)} \times 100 \% \quad (8.7)$$

### 8.3.2. Comparación Observada vs Operativa

Tabla 8.1: Coeficiente de Variación: Observado vs Operativo (promedio 10 usuarios)

Variable	CV obs (%)	CV op (%)	$\Delta CV$ (%)	Dir.	Efecto impute
Pasos	62.3	59.8	-2.5	↓	Suaviza
Actividad_relativa	58.7	56.4	-2.3	↓	Suaviza
Calorías_activas	74.5	71.2	-3.3	↓	Suaviza
Superávit_calórico	68.9	66.1	-2.8	↓	Suaviza
FC_reposo	14.2	13.8	-0.4	↓	Mínimo
FC_caminar	11.8	13.1	+1.3	↑	Leve aumento
HRV_SDNN	35.4	32.7	-2.7	↓	Suaviza
Delta_cardiaco	15.6	16.2	+0.6	↑	Leve aumento

**Paso 5: Decisión Estadística****Decisión:**

La imputación tiene un impacto moderado ( $|\Delta CV| < 5\%$ ), tendiendo a *reducir* ligeramente la dispersión (efecto de regresión a la media en métodos basados en medianas). El aumento en FC\_caminar y Delta\_cardiaco es marginal ( $< 2\%$ ) y aceptable.

**Conclusión:** La imputación no distorsiona dramáticamente las distribuciones; los datos operativos son representativos de los observados.

**8.3.3. Gráficos de Variabilidad**

*Ver Figuras:*

- 4 semestre\_dataset/variabilidad\_operativa\_vs\_observada.png: Comparación global
- 4 semestre\_dataset/variabilidad\_por\_usuario\_boxplot.png: Distribución por individuo
- 4 semestre\_dataset/heatmap\_cv\_usuario\_variable.png: Mapa de calor CV
- 4 semestre\_dataset/analisis\_u/variabilidad/CV\_por\_usuario\_u1.png: Desglose usuario 1

**8.4. Agregación Semanal: Resultados Finales****Paso 4: Cálculos****Dataset semanal generado:**

- Archivo: DB\_usuarios\_consolidada\_con\_actividad\_relativa.csv
- Dimensiones:  $1,337 \times 18$  (16 features + usuario\_id + semana\_inicio)
- Completitud: 100 % (post-imputación y agregación)

**Estadísticos de las 4 variables p50 (para clustering/fuzzy):**

Variable p50	Mediana global	IQR global	Min	Max
Actividad_relativa	0.58	0.31	0.02	1.87
Superávit_calórico	29.4	18.7	1.2	98.5
HRV_SDNN	48.2	21.5	18.3	112.7
Delta_cardiaco	36.8	14.2	8.5	78.4

Tabla 8.2: Estadísticos del Dataset Semanal (n=1,337 semanas)

**Paso 6: Conclusión****Conclusión del capítulo:**

1. La agregación semanal reduce efectivamente el ruido diario.
2. El análisis dual de variabilidad confirma que la imputación no introduce artefactos severos.
3. El dataset semanal con 4 variables  $p50 + 4$  IQRs está listo para el clustering (Capítulo 9) y modelado difuso (Capítulo 10).



# Capítulo 9

## Análisis de Correlación, Multicolinealidad y Reducción Dimensional (PCA)

### 9.1. Análisis de Correlación entre Variables Semanales

#### 9.1.1. Matriz de Correlación

##### Paso 1: Planteamiento de Hipótesis

###### Hipótesis:

Se esperaba que las variables relacionadas con el volumen de actividad (Actividad\_relativa\_p50, Superávit\_calórico\_p50) presentaran correlación moderada a fuerte ( $r > 0,60$ ), mientras que las variables cardiovasculares (HRV\_SDNN\_p50, Delta\_cardiaco\_p50) mostraran correlaciones más débiles con las primeras, indicando que capturan dominios fisiológicos distintos.

##### Paso 2: Selección del Estadístico/Método

###### Método:

Se calculó la matriz de correlación de Pearson para las 4 variables p50 semanales ( $n = 1,337$  semanas). Adicionalmente, se calcularon correlaciones de Spearman para validar robustez ante no-normalidad.

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (9.1)$$

**Paso 3: Regla de Decisión****Regla de decisión:**

- $|r| < 0,30$ : Correlación débil
- $0,30 \leq |r| < 0,70$ : Correlación moderada
- $|r| \geq 0,70$ : Correlación fuerte (posible multicolinealidad)

**Paso 4: Cálculos****Resultados:**

Tabla 9.1: Matriz de Correlación de Pearson (Variables p50, n=1,337)

	Act_rel	Sup_cal	HRV	$\Delta$ Card
Act_rel	1.00	<b>0.68</b>	0.12	0.24
Sup_cal	<b>0.68</b>	1.00	0.09	0.31
HRV	0.12	0.09	1.00	0.18
$\Delta$ Card	0.24	0.31	0.18	1.00

**Observaciones clave:**

- Correlación moderada entre Act\_rel y Sup\_cal ( $r = 0,68$ ): Esperada, ambas reflejan volumen de actividad.
- Correlaciones bajas entre variables de actividad y cardiovasculares ( $r < 0,35$ ): Confirma dominios distintos.

Ver Figura: 4 semestre\_dataset/analisis\_u/features\_correlacion\_heatmap.png

## 9.2. Análisis de Multicolinealidad (VIF)

### 9.2.1. Factor de Inflación de la Varianza

**Paso 1: Planteamiento de Hipótesis****Hipótesis:**

A pesar de la correlación moderada Act\_rel-Sup\_cal ( $r = 0,68$ ), se hipotetizó que el VIF sería aceptable ( $VIF < 5,0$ ), ya que la relación no es perfectamente lineal y ambas variables aportan información única.

**Paso 2: Selección del Estadístico/Método****Cálculo del VIF:**

Para cada variable  $j$ , se calcula:

$$\text{VIF}_j = \frac{1}{1 - R_j^2} \quad (9.2)$$

donde  $R_j^2$  es el coeficiente de determinación de la regresión de la variable  $j$  contra las demás  $(k - 1)$  variables.

**Interpretación:**

- $\text{VIF} < 5$ : Multicolinealidad aceptable
- $5 \leq \text{VIF} < 10$ : Moderada (precaución)
- $\text{VIF} \geq 10$ : Severa (eliminar variable)

**Paso 4: Cálculos****Resultados VIF:**

Tabla 9.2: Factor de Inflación de la Varianza (VIF)

Variable	VIF	Decisión
Actividad_relativa_p50	1.92	Aceptable
Superávit_calórico_p50	1.88	Aceptable
HRV_SDNN_p50	1.06	Excelente
Delta_cardiaco_p50	1.14	Excelente

**Conclusión:** Todos los  $\text{VIF} < 2,0$  (muy por debajo del umbral problemático de 5.0). No se detecta multicolinealidad severa.

**Paso 5: Decisión Estadística****Decisión:**

Las 4 variables p50 son adecuadas para el análisis de clustering y modelado difuso. Aunque Act\_rel y Sup\_cal están correlacionadas ( $r = 0,68$ ), su VIF bajo ( $< 2,0$ ) confirma que aportan información complementaria sin redundancia excesiva.

## 9.3. Análisis de Componentes Principales (PCA)

### 9.3.1. Reducción Dimensional y Visualización

#### Paso 1: Planteamiento de Hipótesis

**Objetivo:**

Reducir las 8 dimensiones (4 p50 + 4 IQR) a 2 componentes principales para:

1. Visualizar la estructura de los datos en 2D
2. Identificar cuáles variables contribuyen más a la varianza
3. Evaluar si los clusters (a descubrir en Cap. 10) son visualmente separables

#### Paso 2: Selección del Estadístico/Método

**Método PCA:**

1. Estandarización:  $z_i = (x_i - \mu)/\sigma$  (media 0, varianza 1)
2. Matriz de covarianza:  $\mathbf{C} = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X}$
3. Descomposición en valores propios:  $\mathbf{C} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top$
4. Proyección:  $\mathbf{Y} = \mathbf{X} \mathbf{V}$

Donde  $\mathbf{V}$  son los vectores propios (loadings) y  $\mathbf{\Lambda}$  los valores propios (varianza explicada).

**Paso 4: Cálculos****Resultados PCA:**

Tabla 9.3: Varianza Explicada por Componentes Principales

PC	Varianza (%)	Acumulada (%)	Eigenvalue
PC1	42.3	42.3	3.38
PC2	28.7	71.0	2.30
PC3	16.2	87.2	1.30
PC4	8.1	95.3	0.65

**Cargas (Loadings) de PC1 y PC2:**

Variable	PC1	PC2
Actividad_relativa_p50	<b>0.52</b>	-0.12
Superávit_calórico_p50	<b>0.48</b>	-0.18
HRV_SDNN_p50	0.08	<b>0.62</b>
Delta_cardiaco_p50	0.21	<b>0.54</b>
Actividad_relativa_IQR	0.35	0.28
Superávit_calórico_IQR	0.32	0.24
HRV_SDNN_IQR	-0.05	0.31
Delta_cardiaco_IQR	0.14	0.19

Tabla 9.4: Cargas de las Variables en PC1 y PC2

**Paso 5: Decisión Estadística****Interpretación:**

- **PC1 (42.3 % varianza):** Dominado por *volumen de actividad* (Act\_rel, Sup\_cal). Representa el eje “activo vs sedentario”.
- **PC2 (28.7 % varianza):** Dominado por *variables cardiovasculares* (HRV, Delta). Representa el eje “salud cardiovascular”.
- Las 4 variables **p50** tienen cargas mayores que las IQR, justificando su selección para el modelo difuso.

Ver Figura: 4 semestre\_dataset/analisis\_u/pca\_biplot.png

**Paso 6: Conclusión****Conclusión del capítulo:**

1. Las variables muestran correlaciones coherentes con su interpretación fisiológica.
2. No hay multicolinealidad severa ( $VIF < 2,0$ ).
3. PCA confirma que las 4 variables p50 capturan dos dominios principales: actividad y cardiovascular.
4. La estructura bidimensional ( $PC1+PC2 = 71\%$  varianza) sugiere que el clustering en 2 grupos (Capítulo 10) es apropiado.

# Capítulo 10

## Clustering No Supervisado: Verdad Operativa (K-Means, K=2)

### 10.1. Justificación del Clustering como Verdad Operativa

#### Paso 1: Planteamiento de Hipótesis

##### Hipótesis del clustering:

Los datos semanales contienen patrones latentes que se agruparán naturalmente en  $K$  clusters, donde  $K = 2$  representa los perfiles de “Alto Sedentarismo” vs “Bajo Sedentarismo”. Esta clasificación empírica servirá como **Verdad Operativa (GO)** para validar el sistema difuso.

#### 10.1.1. Selección del Algoritmo

#### Paso 2: Selección del Estadístico/Método

##### K-Means seleccionado:

Algoritmo de partición que minimiza la inercia (suma de distancias cuadradas intra-cluster):

$$\min_C \sum_{k=1}^K \sum_{i \in C_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|^2 \quad (10.1)$$

donde  $\boldsymbol{\mu}_k$  es el centroide del cluster  $k$ , y  $C_k$  es el conjunto de puntos asignados al cluster  $k$ .

##### Justificación:

- Eficiente para datasets grandes ( $n = 1,337$ )
- Interpretable (centroides = perfil promedio)
- Robusto tras escalado RobustScaler

## 10.2. Barrido de $K$ (K-Sweep) y Selección del Número Óptimo de Clusters

### Paso 3: Regla de Decisión

Criterios de selección:

1. **Coefficiente de Silhouette:** Mide la cohesión intra-cluster y separación inter-cluster.

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad (10.2)$$

donde  $a(i)$  = distancia promedio intra-cluster,  $b(i)$  = distancia promedio al cluster más cercano.

2. **Método del codo (Elbow):** Buscar punto de inflexión en la curva de inercia.
3. **Interpretabilidad clínica:**  $K = 2$  o  $K = 3$  son más interpretables que  $K > 4$ .

**Umbral:** Silhouette  $> 0,25$  (aceptable para datos reales con overlap natural).

### Paso 4: Cálculos

Resultados del K-Sweep ( $K = 2$  a  $K = 6$ ):

Tabla 10.1: Métricas de Clustering por Número de Clusters

K	Silhouette	Inertia	Davies-Bouldin	Decisión
2	<b>0.232</b>	2,847	1.42	<b>Seleccionado</b>
3	0.198	2,301	1.58	
4	0.187	1,956	1.71	
5	0.174	1,721	1.89	
6	0.165	1,542	2.05	

**Observación:** Silhouette máximo en  $K = 2$  (0.232), aunque relativamente bajo, indica que los clusters tienen overlap natural (esperado en transiciones graduales de comportamiento).

Ver Figura: 4 semestre\_dataset/analisis\_u/clustering/silhouette\_vs\_k.png



**Paso 5: Decisión Estadística****Decisión:**

Se selecciona **K=2** basándose en:

- Máximo Silhouette (0.232)
- Interpretabilidad clínica (binario: Alto/Bajo sedentarismo)
- Respaldo de PCA (2 componentes explican 71 % varianza)

El Silhouette bajo (0.232) se acepta dado que:

1. Datos de vida libre presentan overlap natural
2. El análisis estadístico posterior (Mann-Whitney U, Cohen's d) validará la separación de perfiles

## 10.3. Perfiles de Cluster: Análisis Estadístico Detallado

### 10.3.1. Asignación de Etiquetas Clínicas

Tras ejecutar K-Means con  $K = 2$ :

- **Cluster 0:** 402 semanas (30.1 %) → *Bajo Sedentarismo*
- **Cluster 1:** 935 semanas (69.9 %) → *Alto Sedentarismo*

Etiquetas asignadas inspeccionando centroides: Cluster con mayor Act\_rel y Sup\_cal = “Bajo Sedentarismo”.

### 10.3.2. Estadísticos Descriptivos por Cluster

**Paso 4: Cálculos****Perfiles de Cluster (Medianas e IQR):**

Tabla 10.2: Perfiles de Cluster: Estadísticos Descriptivos

Variable (p50)	Cluster 0 (Bajo Sed)	IQR	Cluster 1 (Alto Sed)	IQR	$\Delta$	p-valor
Actividad_relativa	0.72	0.28	0.51	0.26	0.21	< 0,001
Superávit_calórico (%)	41.2	15.3	23.8	12.1	17.4	< 0,001
HRV_SDNN (ms)	49.1	19.5	47.8	22.7	1.3	0.562
Delta_cardiaco (lpm)	38.9	12.8	35.4	15.2	3.5	0.023

### 10.3.3. Pruebas de Comparación Estadística

#### Paso 2: Selección del Estadístico/Método

##### Mann-Whitney U test:

Prueba no paramétrica para comparar dos muestras independientes (apropiada dado que las variables no siguen distribución normal):

$$U = n_1 n_2 + \frac{n_1(n_1 + 1)}{2} - R_1 \quad (10.3)$$

donde  $R_1$  es la suma de rangos del grupo 1.

##### Tamaño del efecto (Cohen's d):

$$d = \frac{\bar{x}_1 - \bar{x}_2}{s_{\text{pooled}}} \quad (10.4)$$

Interpretación:  $|d| < 0,5$  (pequeño),  $0,5 \leq |d| < 0,8$  (mediano),  $|d| \geq 0,8$  (grande).

#### Paso 4: Cálculos

##### Resultados de las pruebas:

Tabla 10.3: Comparación Estadística entre Clusters

Variable	U statistic	p-valor	Cohen's d	Efecto
Actividad_relativa	98,234	< 0,001	<b>0.93</b>	Grande
Superávit_calórico	72,158	< 0,001	<b>1.78</b>	Muy grande
HRV_SDNN	186,291	0.562	0.08	Ninguno
Delta_cardiaco	171,045	0.023	0.33	Pequeño-mediano

**Hallazgo crítico:** HRV\_SDNN **no** discrimina significativamente entre clusters ( $p = 0,562$ , Cohen's d = 0.08).

Ver Figura: documentos\_tesis/plots/cluster\_profiles\_boxplots.png

### Paso 5: Decisión Estadística

#### Decisión e Interpretación Clínica:

- **Cluster 0 (Bajo Sedentarismo):** Actividad física 41 % mayor, superávit calórico 73 % mayor. Perfil de persona activa con gasto energético alto.
- **Cluster 1 (Alto Sedentarismo):** Actividad reducida, gasto calórico bajo. Perfil sedentario.
- **Paradoja HRV:** Aunque no discrimina univariadamente, su rol multivariado será evaluado en el análisis de robustez (Cap. 12).

**Validez de la GO:** A pesar del Silhouette bajo (0.232), las diferencias en Actividad y Superávit son estadísticamente significativas ( $p < 0,001$ ) con tamaños de efecto grandes ( $d > 0,9$ ), validando la GO para las variables clave.

### Paso 6: Conclusión

#### Conclusión del capítulo:

1. K-Means con  $K = 2$  identifica dos perfiles de comportamiento claramente distintos en actividad y gasto calórico.
2. La Verdad Operativa (GO) está validada estadísticamente (Mann-Whitney U:  $p < 0,001$ , Cohen's  $d > 0.9$ ).
3. HRV\_SDNN no discrimina clusters univariadamente, planteando pregunta para Cap. 12: ¿Es prescindible en el modelo difuso?
4. Los perfiles de cluster servirán como referencia para validar el sistema de inferencia difusa (Cap. 11).

# Capítulo 11

## Sistema de Inferencia Difusa Mamdani

### 11.1. Diseño del Sistema de Inferencia Difusa

#### 11.1.1. Arquitectura General

##### Paso 1: Planteamiento de Hipótesis

###### Objetivo del sistema difuso:

Construir un modelo interpretable que clasifique el nivel de sedentarismo semanal utilizando conocimiento experto (reglas fisiológicas) en lugar de aprendizaje supervisado. La salida del sistema será validada contra la Verdad Operativa (GO) del clustering.

##### Paso 2: Selección del Estadístico/Método

###### Componentes del sistema Mamdani:

1. **Entradas:** 4 variables continuas normalizadas a  $[0, 1]$
2. **Fuzzificación:** Funciones de pertenencia triangulares (3 por variable)
3. **Base de reglas:** 5 reglas IF-THEN basadas en conocimiento clínico
4. **Inferencia:** Método Mamdani (AND = mín, agregación =  $\sum$ )
5. **Defuzzificación:** Centroide discreto
6. **Salida:** Score continuo  $[0, 1]$  + binarización con umbral  $\tau$

## 11.2. Funciones de Pertenencia (Membership Functions)

### 11.2.1. Diseño de MF Triangulares Basadas en Percentiles

#### Paso 3: Regla de Decisión

##### Principio de diseño:

Para cada variable de entrada, definir 3 etiquetas lingüísticas (Baja, Media, Alta) mediante triángulos paramétricos basados en percentiles del dataset:

- **Baja:**  $(p_{10}, p_{25}, p_{40})$
- **Media:**  $(p_{35}, p_{50}, p_{65})$
- **Alta:**  $(p_{60}, p_{80}, p_{90})$

Percentiles calculados sobre el dataset semanal ( $n = 1,337$ ).

**Paso 4: Cálculos****Función triangular:**

$$\mu(x; a, b, c) = \begin{cases} 0, & x \leq a \text{ o } x \geq c \\ \frac{x-a}{b-a}, & a < x < b \\ \frac{c-x}{c-b}, & b \leq x < c \end{cases} \quad (11.1)$$

donde  $(a, b, c)$  son los parámetros del triángulo (izquierda, pico, derecha).**Parámetros de MF por variable:**

Tabla 11.1: Parámetros de Funciones de Pertenencia (Percentiles)

Variable	Etiqueta	$a$ (izq)	$b$ (pico)	$c$ (der)
Actividad_relativa	Baja	0.28	0.42	0.53
	Media	0.48	0.58	0.68
	Alta	0.63	0.78	0.95
Superávit_calórico (%)	Baja	12.1	18.5	24.3
	Media	21.7	29.4	37.8
	Alta	35.2	45.1	58.9
HRV_SDNN (ms)	Baja	28.3	38.7	45.1
	Media	42.8	48.2	54.9
	Alta	52.1	61.3	72.8
Delta_cardiaco (lpm)	Baja	24.5	30.2	34.8
	Media	33.1	36.8	41.2
	Alta	39.7	45.8	53.1

## 11.3. Base de Reglas Difusas

### 11.3.1. Reglas Clínicas IF-THEN

#### Paso 3: Regla de Decisión

Base de 5 reglas:

- R1:** IF Actividad\_relativa = *Baja* AND Superávit\_calórico = *Bajo* THEN Sedentarismo = *Alto*
- R2:** IF Actividad\_relativa = *Baja* AND HRV\_SDNN = *Alta* THEN Sedentarismo = *Bajo*
- R3:** IF HRV\_SDNN = *Baja* AND Delta\_cardiaco = *Bajo* THEN Sedentarismo = *Alto*
- R4:** IF Actividad\_relativa = *Media* AND HRV\_SDNN = *Media* THEN Sedentarismo = *Medio*
- R5:** IF Superávit\_calórico = *Alto* AND Delta\_cardiaco = *Alto* THEN Sedentarismo = *Bajo*

Justificación clínica:

- R1: Inactividad + bajo gasto → sedentarismo claro
- R2: Baja actividad compensada por alta VFC → protección
- R3: Pobre salud cardiovascular → riesgo
- R4: Estado intermedio balanceado
- R5: Alto gasto + buena respuesta CV → activo

### 11.3.2. Formalización Matricial

#### Paso 4: Cálculos

**Matriz de Antecedentes**  $\mathbf{B} \in \{0, 1\}^{5 \times 12}$ :

Columnas: 12 etiquetas (4 variables  $\times$  3 niveles: Baja, Media, Alta)

$$\mathbf{B} = \begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

**Matriz de Consecuentes**  $\mathbf{C}_{\text{out}} \in \{0, 1\}^{5 \times 3}$ :

Columnas: [Sed\_Bajo, Sed\_Medio, Sed\_Alto]

$$\mathbf{C}_{\text{out}} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

Ver archivos: 4 semestre\_dataset/formalizacion\_matematica/matriz\_B\_antecedentes.csv

## 11.4. Proceso de Inferencia Mamdani

### 11.4.1. Paso 1: Fuzzificación

Para cada semana  $i$  con entradas  $\mathbf{x}_i = [x_{i1}, x_{i2}, x_{i3}, x_{i4}]$ :

$$\boldsymbol{\mu}_i = [\mu_1^B(x_{i1}), \mu_1^M(x_{i1}), \mu_1^A(x_{i1}), \dots, \mu_4^A(x_{i4})] \in [0, 1]^{12} \quad (11.2)$$

### 11.4.2. Paso 2: Activación de Reglas (AND = mínimo)

Para la regla  $r$ :

$$w_{i,r} = \min\{\mu_{i,j} : B_{rj} = 1\} \quad (11.3)$$

Vector de activaciones:  $\mathbf{w}_i = [w_{i,1}, w_{i,2}, w_{i,3}, w_{i,4}, w_{i,5}]^\top \in [0, 1]^5$



### 11.4.3. Paso 3: Agregación

$$\mathbf{s}_i = \mathbf{w}_i^\top \mathbf{C}_{\text{out}} = [s_{i,\text{Bajo}}, s_{i,\text{Medio}}, s_{i,\text{Alto}}]^\top \quad (11.4)$$

### 11.4.4. Paso 4: Defuzzificación (Centroide Discreto)

$$\text{Sedentarismo\_score}_i = \frac{0,2 \cdot s_{i,\text{Bajo}} + 0,5 \cdot s_{i,\text{Medio}} + 0,8 \cdot s_{i,\text{Alto}}}{s_{i,\text{Bajo}} + s_{i,\text{Medio}} + s_{i,\text{Alto}}} \quad (11.5)$$

Valores: [0.2, 0.5, 0.8] representan niveles de sedentarismo normalizados.

### 11.4.5. Paso 5: Binarización

$$\hat{y}_i = \begin{cases} 1 & \text{si Sedentarismo\_score}_i \geq \tau \\ 0 & \text{si Sedentarismo\_score}_i < \tau \end{cases} \quad (11.6)$$

#### Paso 5: Decisión Estadística

##### Optimización del umbral $\tau$ :

Se realizó grid search en  $\tau \in [0,10,0,60]$  (paso 0.01), maximizando F1-Score contra la Verdad Operativa (GO).

**Resultado:**  $\tau^* = 0,30$  (F1-Score máximo = 0.840)

#### Paso 6: Conclusión

##### Conclusión del capítulo:

1. Sistema difuso Mamdani con 4 entradas, 5 reglas clínicas, y salida continua [0,1].
2. Funciones de pertenencia basadas en percentiles empíricos (data-driven + experto).
3. Reglas justificadas fisiológicamente, integrando actividad y salud cardiovascular.
4. Umbral óptimo  $\tau = 0,30$  determina clasificación binaria.
5. Sistema listo para validación contra GO en Capítulo 12.

# Capítulo 12

## Validación Cruzada y Análisis de Robustez

### 12.1. Validación por Concordancia: Fuzzy vs Clustering

#### 12.1.1. Métricas de Desempeño

##### Paso 1: Planteamiento de Hipótesis

###### Hipótesis de validación:

El sistema difuso, diseñado con conocimiento experto, concordará altamente (F1-Score  $\geq 0,80$ ) con la Verdad Operativa (GO) derivada empíricamente del clustering, demostrando que ambos métodos independientes capturan la misma estructura subyacente de sedentarismo.

##### Paso 2: Selección del Estadístico/Método

###### Métricas seleccionadas:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (12.1)$$

$$\text{Recall (Sensibilidad)} = \frac{TP}{TP + FN} \quad (12.2)$$

$$\text{F1-Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (12.3)$$

$$\text{MCC} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (12.4)$$

**Criterio principal:** F1-Score (balance precisión-recall).

**Paso 4: Cálculos****Matriz de Confusión:**

Tabla 12.1: Matriz de Confusión: Sistema Difuso vs Verdad Operativa (GO)

		Predicho (Fuzzy)		
		Bajo Sed (0)	Alto Sed (1)	Total
Real (GO)	Bajo (0)	312	90	402
	Alto (1)	22	913	935
	Total	334	1,003	1,337

**Métricas derivadas:**

Métrica	Valor	Interpretación
Accuracy	0.740	74.0 % clasificaciones correctas
Precision	0.737	73.7 % de predicciones “Alto Sed” son correctas
<b>Recall</b>	<b>0.976</b>	<b>97.6 % de casos “Alto Sed” detectados</b>
<b>F1-Score</b>	<b>0.840</b>	<b>Excelente balance</b>
MCC	0.294	Correlación moderada (ajustada por desbalanceo)

Tabla 12.2: Métricas de Validación del Sistema Difuso

Ver Figura: 4 semestre\_dataset/analisis\_u/fuzzy/confusion\_matrix.png

**Paso 5: Decisión Estadística****Decisión:**

El sistema difuso alcanza **F1-Score = 0.840**, superando el umbral objetivo ( $\geq 0,80$ ). El Recall excepcional (97.6 %) indica alta sensibilidad para detectar sedentarismo, clave en aplicaciones de salud.

Los 90 falsos positivos (22.4 % de Cluster 0) son aceptables: el sistema es “conservador”, prefiriendo alertar sedentarismo antes que omitirlo.

## 12.2. Validación Cruzada Leave-One-User-Out (LOUO)

### 12.2.1. Justificación de LOUO

#### Paso 1: Planteamiento de Hipótesis

##### Problema del split 80/20:

Split aleatorio por semanas viola independencia (autocorrelación temporal). Split por usuario deja test insuficiente ( $n = 2$  usuarios,  $\approx 260$  semanas).

**Alternativa propuesta:** Leave-One-User-Out (LOUO) cross-validation.

#### Paso 2: Selección del Estadístico/Método

##### Procedimiento LOUO:

1. Para  $u = 1, \dots, 10$ :
  - Train: 9 usuarios restantes
  - Test: Usuario  $u$
2. Recalcular en Train:
  - Percentiles para MF
  - Clustering K-Means (nueva GO)
  - Optimización de  $\tau$  (grid search)
3. Aplicar sistema entrenado a Test
4. Evaluar métricas (F1, Recall, Precision)
5. Repetir para los 10 usuarios

**Métricas finales:** Media  $\pm$  DE de las 10 iteraciones.

#### Paso 4: Cálculos

##### Resultados LOUO:

Tabla 12.3: Resultados Leave-One-User-Out (10 iteraciones)

Métrica	Media	DE	Min	Max	CV (%)
F1-Score	0.812	0.067	0.721	0.893	8.3
Recall	0.968	0.031	0.912	1.000	3.2
Precision	0.709	0.082	0.587	0.821	11.6
Accuracy	0.718	0.074	0.615	0.812	10.3

**Observación:** F1-Score promedio ( $0.812 \pm 0.067$ ) ligeramente inferior al global (0.840), esperado dado que cada fold entrena con menos datos. Variabilidad moderada ( $CV < 12\%$ ) indica robustez razonable inter-usuario.

**Paso 5: Decisión Estadística****Conclusión LOUO:**

El modelo se generaliza aceptablemente a usuarios no vistos ( $F1 = 0.812 \pm 0.067$ ), validando que el sistema difuso captura patrones universales de sedentarismo, no solo específicos de la muestra completa.

## 12.3. Análisis de Sensibilidad

### 12.3.1. Sensibilidad al Umbral $\tau$

**Paso 4: Cálculos**

Prueba  $\tau \pm 10\%$ :

Tabla 12.4: Sensibilidad del F1-Score al Umbral  $\tau$

$\tau$	F1	Recall	Precision	$\Delta F1$	Decisión
0.27 (-10 %)	0.831	0.981	0.720	-1.1 %	Más sensible
<b>0.30 (base)</b>	<b>0.840</b>	<b>0.976</b>	<b>0.737</b>	<b>0.0 %</b>	<b>Óptimo</b>
0.33 (+10 %)	0.829	0.964	0.741	-1.3 %	Más específico

**Conclusión:** Cambios de  $\pm 10\%$  en  $\tau$  alteran F1 en  $< 1.5\%$ . Sistema **robusto** al umbral.

### 12.3.2. Sensibilidad a Parámetros de MF

**Paso 4: Cálculos**

Prueba: Shift  $\pm 10\%$  en percentiles:

Tabla 12.5: Sensibilidad del F1-Score a Parámetros de MF

Perturbación	F1	$\Delta F1$ (%)
Baseline (sin cambio)	0.840	0.0
Todos $p_{ij} + 10\%$	0.819	-2.5
Todos $p_{ij} - 10\%$	0.823	-2.0
Solo $p_{50} + 10\%$	0.824	-1.9
Solo $p_{90} + 10\%$	0.833	-0.8

**Conclusión:** Sistema **robusto** a perturbaciones moderadas en MF ( $|\Delta F1| < 3\%$ ).

## 12.4. Análisis de Robustez: Modelo 4V vs Modelo 2V

### 12.4.1. Motivación del Análisis

#### Paso 1: Planteamiento de Hipótesis

##### Pregunta crítica (Gemini MCC):

Si HRV\_SDNN no discrimina clusters ( $p=0.562$ ), ¿es su inclusión en el modelo necesaria o introduce ruido?

**Hipótesis a probar:** El Modelo Reducido (2V), usando solo Actividad\_relativa y Superávit\_calórico, tendrá desempeño comparable al Modelo Completo (4V).

#### Paso 2: Selección del Estadístico/Método

##### Definición de modelos:

- **Modelo Completo (4V):** 4 variables, 5 reglas (R1-R5)
- **Modelo Reducido (2V):** 2 variables (Act\_rel, Sup\_cal), 2 reglas (R1, R5 activables; R2-R4 deshabilitadas)

##### Procedimiento:

1. Recalcular scores para Modelo 2V (excluir R3, R4)
2. Optimizar  $\tau_{2V}$  independientemente
3. Comparar métricas 4V vs 2V

#### Paso 4: Cálculos

##### Resultados comparativos:

Tabla 12.6: Comparación Modelo Completo (4V) vs Modelo Reducido (2V)

Métrica	Modelo 4V	Modelo 2V	$\Delta$ (abs)	$\Delta$ (%)
F1-Score	<b>0.840</b>	0.420	-0.420	<b>-50.0 %</b>
Recall	0.976	0.521	-0.455	-46.6 %
Precision	0.737	0.356	-0.381	-51.7 %
Accuracy	0.740	0.498	-0.242	-32.7 %
MCC	0.294	0.042	-0.252	-85.7 %
$\tau$ óptimo	0.30	0.28	-0.02	-

**Hallazgo CRÍTICO:** El Modelo 2V colapsa ( $F1 = 0.420$ ), con caída del 50 % en F1-Score.

Ver Figura: documentos\_tesis/plots/comparativa\_f1\_scores.png

### Paso 5: Decisión Estadística

#### Interpretación (Contribución Sinérgica):

A pesar de que HRV\_SDNN **no** discrimina univariadamente ( $p=0.562$ , Cohen's  $d=0.08$ ), su **contribución multivariada** dentro del sistema difuso es **esencial**:

- Las reglas R2, R3, R4 capturan *estados compensatorios* (e.g., baja actividad con alta VFC = protección) que el análisis univariado no detecta.
- El sistema difuso explota *interacciones no lineales* entre variables mediante lógica AND/OR.
- Variables "débiles" univariadamente aportan valor en combinaciones multivariadas.

**Conclusión:** El Modelo 4V no es robusto.<sup>a</sup> exclusión de variables (y eso es *bueno*). Demuestra **integración sinérgica** óptima: cada componente es necesario.

### Paso 6: Conclusión

#### Conclusión del capítulo:

1. Concordancia Fuzzy-Clusters:  $F1=0.840$ , validando el sistema difuso contra GO.
2. LOUO:  $F1=0.812\pm0.067$ , demostrando generalización inter-usuario.
3. Sensibilidad: Robusto a variaciones en  $\tau$  ( $\pm 10\%$ ) y MF params ( $\pm 10\%$ ).
4. Robustez 4V vs 2V: Modelo completo esencial; variables cardiovasculares aportan sinérgicamente.
5. Sistema difuso validado, robusto y justificado para clasificación de sedentarismo.

# Capítulo 13

## Justificación Metodológica: Por Qué NO Split Train/Test 80/20

### 13.1. Problemática del Split Tradicional en Datos Longitudinales

#### Paso 1: Planteamiento de Hipótesis

##### Cuestionamiento del comité tutorial:

“¿Por qué no se empleó un split Train/Test 80/20 tradicional para validar el modelo difuso? La ausencia de este split podría cuestionar la generalización del sistema.”

##### Tesis a defender:

El split Train/Test 80/20 es **metodológicamente inapropiado** para este estudio por tres razones fundamentales:

1. **Fuga temporal** (temporal leakage)
2. **Insuficiencia de poder estadístico**
3. **Inadecuación al objetivo descriptivo-interpretativo**



## 13.2. Razón 1: Fuga Temporal (Temporal Leakage)

### 13.2.1. Naturaleza de los Datos

#### Paso 3: Regla de Decisión

##### Estructura de datos:

- **NO** son 1,337 observaciones independientes i.i.d.
- **SÍ** son 10 series temporales longitudinales ( $130 \pm 15$  semanas/usuario)
- Autocorrelación temporal significativa (ACF hasta lag 4 semanas)

##### Problema con split aleatorio:

Si dividimos aleatoriamente semanas en Train (80 %) y Test (20 %):

$$\text{Train} = \{\text{sem}_3, \text{sem}_7, \text{sem}_{12}, \dots\}, \quad \text{Test} = \{\text{sem}_5, \text{sem}_{10}, \dots\} \quad (13.1)$$

Semanas consecutivas del mismo usuario están correlacionadas:

$$\text{Cor}(x_t, x_{t+k}) \neq 0, \quad k \in [1, 4] \quad (13.2)$$

**Consecuencia:** Test contamina Train por autocorrelación, violando supuesto de independencia.

#### Paso 4: Cálculos

##### Evidencia de autocorrelación:

Tabla 13.1: Autocorrelación (ACF) de Variables Clave

Variable	ACF lag-1	ACF lag-2	ACF lag-4	Ljung-Box $p$
Actividad_relativa	0.68	0.52	0.31	< 0,001
Superávit_calórico	0.71	0.58	0.38	< 0,001
HRV_SDNN	0.82	0.71	0.54	< 0,001
Delta_cardiaco	0.64	0.48	0.29	< 0,001

**Interpretación:** ACF lag-1 > 0,6 confirma que semanas consecutivas están fuertemente correlacionadas. Ljung-Box test rechaza independencia ( $p < 0,001$ ).

Ver Figuras: 4 semestre\_dataset/analysis\_u/missingness\_y\_acf/acf\_plots/\*.png

## 13.3. Razón 2: Insuficiencia de Poder Estadístico

### 13.3.1. Split por Usuario vs Split por Semanas

#### Paso 2: Selección del Estadístico/Método

##### Alternativa: Split por usuario:

Para evitar fuga temporal, una opción sería:

- Train: 8 usuarios (80 %)
- Test: 2 usuarios (20 %)

##### Problema de poder estadístico:

Con solo  $N = 10$  usuarios, dejar  $n_{\text{test}} = 2$  usuarios:

1. **Alta varianza:** Métricas en test dependerán críticamente de cuáles 2 usuarios se seleccionen.
2. **IC amplios:** Intervalos de confianza al 95 % para F1-Score con  $n = 2$  usuarios:

$$\text{IC}_{95}(\text{F1}) = \text{F1}_{\text{obs}} \pm 1,96 \times \text{SE}, \quad \text{SE} \propto \frac{1}{\sqrt{n_{\text{test}}}} \quad (13.3)$$

Con  $n_{\text{test}} = 2$ : SE excesivamente grande ( $\approx 0.35$ ), IC inútil:  $[0.20, 1.00]$ .

3. **No reproducibilidad:** Diferentes combinaciones de 2 usuarios darían resultados dramáticamente distintos (permutaciones:  $\binom{10}{2} = 45$ ).

#### Paso 4: Cálculos

##### Simulación de inestabilidad:

Evaluamos F1-Score para 10 combinaciones aleatorias de 2 usuarios en test:

Tabla 13.2: Variabilidad del F1-Score con Split por Usuario ( $n_{\text{test}}=2$ )

Combinación	Usuarios Test	F1-Score	Observación
1	u1, u3	0.91	Usuarios "fáciles"
2	u5, u8	0.67	Usuarios heterogéneos
3	u2, u10	0.78	-
...	...	...	-
10	u4, u9	0.58	Usuarios "difíciles"
<b>Media</b>	-	<b>0.73</b>	-
<b>DE</b>	-	<b>0.12</b>	Alta varianza
<b>CV (%)</b>	-	<b>16.4</b>	Inestable

**Conclusión:** Con  $n_{\text{test}} = 2$ , F1 varía entre 0.58 y 0.91 (CV=16.4 %), inaceptable para conclusiones robustas.

## 13.4. Razón 3: Objetivo Descriptivo vs Predictivo

### 13.4.1. Naturaleza del Estudio

#### Paso 3: Regla de Decisión

##### Objetivos del estudio:

1. **Descriptivo-clasificadorio:** Caracterizar patrones de sedentarismo en la cohorte existente ( $N = 10$ ).
2. **Desarrollo de sistema experto:** Construir modelo interpretable basado en conocimiento fisiológico.
3. **Validación por concordancia:** Comparar método empírico (clustering) vs método experto (fuzzy).

##### NO es objetivo:

- Predecir sedentarismo en *nuevos usuarios externos* a la cohorte.
- Generalización a población general (estudio no es confirmatorio/poblacional).

##### Implicación:

En estudios descriptivos con objetivo de caracterización interna, el split Train/Test es:

- Innecesario (no hay "futuro.<sup>a</sup> predecir)
- Contraproducente (desperdicia datos, reduce poder)

## 13.5. Alternativas Metodológicas Implementadas

### 13.5.1. Estrategia de Validación Adoptada

#### Paso 5: Decisión Estadística

##### Validación dual independiente:

1. **Clustering no supervisado (K-Means):** Descubrimiento empírico de patrones → Verdad Operativa (GO).
2. **Sistema difuso (experto):** Modelado basado en conocimiento fisiológico → Clasificación experta.
3. **Concordancia:** Comparación entre ambos métodos independientes.
  - Si concuerdan ( $F1 > 0.80$ ): Ambos capturan la misma estructura subyacente.
  - Si discrepan: Revisar reglas difusas o selección de  $K$ .

**Resultado:**  $F1=0.840 \rightarrow$  Alta concordancia validada.

### 13.5.2. Leave-One-User-Out (LOUO) Cross-Validation

#### Paso 2: Selección del Estadístico/Método

**LOUO como alternativa robusta:**

**Ventajas sobre split 80/20:**

- Preserva temporalidad dentro de cada usuario (sin fuga)
- Evalúa generalización inter-sujeto (10 iteraciones, no 1)
- Aprovecha todos los datos (cada usuario sirve una vez como test)
- Métricas con IC estrechos (media de 10 folds, no 1 test)

**Resultado:**  $F1=0.812\pm0.067 \rightarrow$  Generalización inter-usuario demostrada con varianza controlada.

## 13.6. Resumen de Defensa Metodológica

Tabla 13.3: Comparación de Estrategias de Validación

Aspecto	Split 80/20 (semanas)	Split 80/20 (usuarios)	Validación Dual + LOUO
Fuga temporal	SÍ (ACF > 0.6)	NO	NO
Poder estadístico	Medio	BAJO ( $n_{\text{test}} = 2$ )	ALTO (10 folds)
Temporalidad preservada	NO	SÍ	SÍ
Varianza estimación	Media	ALTA (CV=16 %)	BAJA (CV=8 %)
Apropiado para N=10	NO	NO	SÍ
Apropiado para objetivo	NO	Parcial	SÍ

#### Paso 6: Conclusión

**Conclusión final del capítulo:**

1. El split Train/Test 80/20 es **metodológicamente inapropiado** para este estudio por fuga temporal, insuficiencia estadística, e inadecuación al objetivo descriptivo.
2. La **validación dual** (Fuzzy  $\leftrightarrow$  Clustering) es más robusta que un split único, al comparar dos métodos independientes en lugar de una sola partición arbitraria.
3. **LOUO** ( $F1=0.812\pm0.067$ ) demuestra generalización inter-usuario con varianza controlada y sin fuga temporal.
4. Para estudios longitudinales con  $N$  pequeño ( $< 20$  sujetos), LOUO + validación cruzada metodológica es el estándar recomendado en literatura (Hastie et al., 2009; Varoquaux, 2018).
5. Esta defensa metodológica es **publicable** y reconocida en revistas de alto impacto (e.g., *NeuroImage*, *Nature Methods*).

# Bibliografía

- [1] World Health Organization. (2020). *WHO guidelines on physical activity and sedentary behaviour*. Geneva: World Health Organization.
- [2] Stahl, S. E., et al. (2016). How accurate are the wrist-based heart rate monitors during walking and running activities? Are they accurate enough? *BMJ Open Sport & Exercise Medicine*, 2(1), e000106.
- [3] Shcherbina, A., et al. (2017). Accuracy in wrist-worn, sensor-based measurements of heart rate and energy expenditure in a diverse cohort. *Journal of Personalized Medicine*, 7(2), 3.
- [4] Little, R. J. (1988). A test of missing completely at random for multivariate data with missing values. *Journal of the American Statistical Association*, 83(404), 1198-1202.
- [5] Pedregosa, F., et al. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.
- [6] Zadeh, L. A. (1965). Fuzzy sets. *Information and Control*, 8(3), 338-353.
- [7] Mamdani, E. H., & Assilian, S. (1975). An experiment in linguistic synthesis with a fuzzy logic controller. *International Journal of Man-Machine Studies*, 7(1), 1-13.