

Task 4

Tema: Tidy-data

Profesor: Eduard F. Martínez González

Septiembre 2025

Introducción

El presente taller ha sido diseñado con el propósito de consolidar y aplicar los conocimientos adquiridos en torno a la limpieza y organización de datos, haciendo uso de la librerías de `dplyr` y `janitor`, y de los contenidos desarrollados en la semana 7 del curso.

Para el desarrollo de las actividades propuestas, se pone a disposición una base de datos compuesta por 1.000 registros de pacientes atendidos en una clínica entre los años 2010 y 2025. Es importante señalar que dicha base de datos fue construida a partir de información sintética, generada de manera aleatoria, y que no corresponde a datos reales. Su única finalidad es servir como insumo para los ejercicios académicos del presente taller. La base de datos incluye variables asociadas a características sociodemográficas y clínicas de los pacientes, tales como edad, género, nivel educativo, afiliación a EPS, peso, estatura, fecha de registro y estado nutricional estimado mediante el índice de masa corporal (IMC).

Repositorio

Para resolver este taller, el estudiante deberá **descargar y descomprimir** el repositorio `task-tidy-data.zip`, disponible en el siguiente enlace: [aquí](#). El repositorio contiene la siguiente estructura:

- Un archivo `proj.Rproj`, el cual debe ser utilizado para **iniciar RStudio**. Esto garantizará que el directorio de trabajo se establezca correctamente en la carpeta principal del repositorio.
- Una carpeta `code`, que incluye el **script base** sobre el cual deberá desarrollar la solución del taller. En caso de requerir librerías adicionales, estas deberán ser instaladas y llamadas al inicio del script.
- Una carpeta `input`, que contiene la **base de datos necesaria** para resolver los ejercicios propuestos.

Se enfatiza que el correcto uso del archivo `proj.Rproj` es indispensable para evitar errores

relacionados con rutas o directorios durante el desarrollo del taller.

Instrucciones

- No seguir correctamente **todas** las instrucciones del taller implicará una penalización del **20 % sobre la nota total**.
- El taller puede desarrollarse en grupos de hasta dos personas. Únicamente uno de los integrantes debe cargar la entrega en la plataforma Intu.
- Este documento presenta dos opciones de trabajo —*Taller 1* y *Taller 2*—. **Cada grupo debe escoger únicamente uno de los dos talleres para su desarrollo.**
- La fecha límite de entrega para el grupo 007 es el **martes XX de septiembre a las 2:00 p.m.** y para el grupo 003 es el **viernes XX de septiembre a las 10:00 a.m.**. A partir de ese momento la plataforma Intu se cerrará y no permitirá subir más archivos.
- La plataforma Intu recibirá exclusivamente dos archivos:
 - Un archivo en formato **.R** (script), con el desarrollo de los puntos del taller, documentado de manera clara y ordenada.
 - Un archivo en formato **PDF** (no Word ni ningún otro formato), que contenga las respuestas interpretativas solicitadas.
- El script en **R** debe cumplir con las siguientes condiciones:
 - Incluir al inicio: los nombres de los integrantes, la versión de R utilizada y la carga de todas las librerías necesarias.
 - Estar limpio y organizado: elimine funciones o líneas de código innecesarias.
 - Documentar claramente el desarrollo de cada punto utilizando comentarios: **## Punto 1, ## Punto 2**, etc.

```
## Nombre(s) de Autor(es)
## R version 4.5.0

## limpiar entorno
rm(list=ls())

## cargar librerías
require(dplyr)
require(skimr)
require(janitor)
require(rio)

## Punto 1
...
## Punto 2
...
```

Opción Taller 1: Análisis e Interpretación de Datos Clínicos

1. (15 pts) Cargue la base de datos y use la función `skim()` para describir las variables. ¿Qué tipo de problemas de calidad se evidencian? (valores faltantes, extremos, inconsistentes).
2. (20 pts) Cree una nueva variable de índice de masa corporal (`imc`) e interprete qué grupo de pacientes (por género o edad) presenta un mayor riesgo nutricional.
3. (20 pts) Realice un análisis descriptivo del número de pacientes por EPS y nivel educativo. ¿Existen categorías mal etiquetadas, vacías o inconsistentes?
4. (15 pts) Analice las fechas de registro. ¿Hay registros antes del año 2010 o posteriores a 2025? ¿Qué posibles errores de digitación o inconsistencias temporales se identifican?
5. (30 pts) Redacte un texto (máximo media página) que resuma los principales retos de limpieza y organización encontrados en el conjunto de datos. Argumente brevemente cómo podrían corregirse.

Total: 100 pts

Opción Taller 2: Transformaciones, Limpieza y Joins con tidyverse

1. (15 pts) Cargue la base de datos y reescriba el siguiente pipeline usando el operador `%>%`:

```
df <- import("input/data-pacientes.csv")
df <- select(df, -nombre)
df <- mutate(df, edad = ifelse(is.na(edad), 0,
                                edad))
```

2. (20 pts) Genere una variable llamada `clasificacion_imc` con base en el valor de `imc`, utilizando `case_when()`. Cree una tabla con el número de pacientes en cada categoría, usando `tabyl()`.
3. (20 pts) Simule una tabla de regiones para ciertos pacientes y combínela con la base principal usando `left_join()`.
4. (15 pts) Cree un pipeline que filtre pacientes mayores de 65 años con IMC mayor a 30, seleccione las variables clave, y exporte el resultado como `'csv'`.

5. (30 pts) Con ayuda de `group_by()`, `summarise()` y `pivot_wider()`, construya una tabla que muestre el promedio de edad y promedio de IMC por género y EPS. Explique qué patrones llaman la atención.

Total: 100 pts