

## Evaluación N°2 – Informe de análisis con R

### Objetivo

Aplicar técnicas de procesamiento y análisis de datos con R, mediante la limpieza, transformación y visualización del conjunto de datos histórico de titulados 2007–2020, con el propósito de desarrollar una comprensión analítica sobre la distribución, evolución y características de los titulados en Chile, a través de la generación de tablas resumen, gráficos estadísticos y conclusiones interpretativas.

### Contexto del archivo

El archivo “titulados\_histórico 2007–2020.csv” contiene información oficial sobre los titulados de educación superior en Chile entre los años 2007 y 2020. Este conjunto de datos recopila información proveniente del Ministerio de Educación, y permite analizar tendencias de titulación por año, género, institución, sede, carrera, área del conocimiento y región, entre otros criterios relevantes.

Cada registro representa una combinación específica de carrera, sede e institución, junto con la cantidad total de titulados (mujeres, hombres y totales), su rango de edad y promedios de edad de titulación.

Este tipo de información es clave para comprender el comportamiento histórico de la educación superior chilena y facilita el desarrollo de análisis descriptivos y visuales mediante herramientas de minería de datos.

### Entrega:

Esta entrega es de tipo individual

**Nombre de archivo:** cada archivo deberá tener su nombre y apellido

- codido\_nombre\_apellido.r
- archivo\_limpio\_nombre\_apellido.csv
- análisis\_nombre\_apellido.pdf

Debe entregar un zip con los 3 documentos, en caso de que cuyo archivo supere el máximo de tamaño permitido por la plataforma debe subir el archivo a **OneDrive(Microsoft)** y compartir link para acceder.

**Fecha de entrega:** miércoles 29 de octubre 23:30 hrs, vía plataforma

**Ponderación:** 30%

**No se recibirán trabajos por correos, si el trabajo no está en la plataforma en la fecha y plazo máximo se le asignara nota mínima.**

## Parte I – Limpieza del archivo (90 pts)

Como base, debe realizar la limpieza del archivo original titulados\_histórico\_2007\_2020.csv.

Tareas obligatorias:

- I. Renombrar columnas (15 pts):

Utilizar los nuevos nombres indicados en la tabla siguiente:

Nombre antiguo	Nuevo nombre
ANNIO	ANNIO
<b>TOTAL TITULADOS</b>	<b>TOTAL_TITULADOS</b>
<b>TITULADOS MUJERES POR PROGRAMA</b>	<b>TITULADOS_MUJERES_POR_PROGRAMA</b>
<b>TITULADOS HOMBRES POR PROGRAMA</b>	<b>TITULADOS_HOMBRES_POR_PROGRAMA</b>
<b>CLASIFICACION INSTITUCION NIVEL 1</b>	<b>CLASIFICACION_NIVEL_1</b>
<b>CLASIFICACION INSTITUCION NIVEL 2</b>	<b>CLASIFICACION_NIVEL_2</b>
<b>CLASIFICACION INSTITUCION NIVEL 3</b>	<b>CLASIFICACION_NIVEL_3</b>
<b>CODIGO INSTITUCION</b>	<b>CODIGO_INSTITUCION</b>
<b>NOMBRE INSTITUCION</b>	<b>NOMBRE_INSTITUCION</b>
<b>COMUNA</b>	<b>COMUNA</b>
<b>PROVINCIA</b>	<b>PROVINCIA</b>
<b>REGION</b>	<b>REGION</b>
<b>NOMBRE SEDE</b>	<b>NOMBRE_SEDE</b>
<b>NOMBRE CARRERA</b>	<b>NOMBRE_CARRERA</b>
<b>AREA DEL CONOCIMIENTO</b>	<b>AREA_CONOCIMIENTO</b>
<b>CINE-F_97 AREA</b>	<b>AREA</b>
<b>CINE-F_97 SUBAREA</b>	<b>SUBAREA</b>
<b>AREA CARRERA GENERICA</b>	<b>AREA_CARRERA_GENERICA</b>
<b>NIVEL GLOBAL</b>	<b>NIVEL_GLOBAL</b>
<b>CARRERA CLASIFICACION NIVEL 1</b>	<b>CARRERA_CLASIFICACION_NIVEL_1</b>
<b>CARRERA CLASIFICACION NIVEL 2</b>	<b>CARRERA_CLASIFICACION_NIVEL_2</b>
<b>MODALIDAD</b>	<b>MODALIDAD</b>
<b>JORNADA</b>	<b>JORNADA</b>
<b>TIPO DE PLAN DE LA CARRERA</b>	<b>TIPO_PLAN_CARRERA</b>
<b>DURACION ESTUDIO CARRERA</b>	<b>DURACION_ESTUDIO_CARRERA</b>
<b>DURACION TOTAL DE LA CARRERA</b>	<b>DURACION_TOTAL_CARRERA</b>
<b>CODIGO CARRERA</b>	<b>CODIGO_CARRERA</b>
<b>TOTAL RANGO EDAD</b>	<b>TOTAL_RANGO_EDAD</b>
<b>RANGO DE EDAD 15 A 19 AÑOS</b>	<b>RANGO_15_A_19_ANNIOS</b>
<b>RANGO DE EDAD 20 A 24 AÑOS</b>	<b>RANGO_20_A_24_ANNIOS</b>
<b>RANGO DE EDAD 25 A 29 AÑOS</b>	<b>RANGO_25_A_29_ANNIOS</b>
<b>RANGO DE EDAD 30 A 34 AÑOS</b>	<b>RANGO_30_A_34_ANNIOS</b>
<b>RANGO DE EDAD 35 A 39 AÑOS</b>	<b>RANGO_35_A_39_ANNIOS</b>

<b>RANGO DE EDAD 40 Y MAS AÑOS</b>	RANGO_40_Y_MAS_ANNIOS
<b>RANGO DE EDAD SIN INFORMACION</b>	RANGO_SIN_INFORMACION
<b>PROMEDIO EDAD CARRERA</b>	PROMEDIO_EDAD_CARRERA
<b>PROMEDIO EDAD MUJER</b>	PROMEDIO_EDAD_MUJER
<b>PROMEDIO EDAD HOMBRE .....</b>	PROMEDIO_EDAD_HOMBRE

- II. Reemplazar valores vacíos, Null y NaN por 0 (15 pts)  
 Asegúrese de convertir correctamente las celdas vacías en valores numéricos válidos.
- III. Corregir registros en PROMEDIO\_EDAD\_HOMBRE (15 pts)  
 Elimine las comas o símbolos sobrantes al final de los registros.
- IV. Modificar columna ANNIO (15 pts)  
 Elimine el prefijo “TIT\_” para dejar únicamente el año numérico.
- V. Eliminar columnas innecesarias (15 pts)
- CODIGO\_INSTITUCION
  - CODIGO\_CARRERA
  - RANGO\_SIN\_INFORMACION
- VI. Guardar DataFrame limpio (15 pts)  
 Guarde el resultado como:  
 archivo\_limpio\_nombre\_apellido.csv

## Parte II – Desarrollo de análisis y visualización (100 pts)

Mediante código en R, deberá responder a las siguientes 10 preguntas, cada una con:

- Una tabla resumen (4 pts)
- Un gráfico asociado (4 pts)
- Un breve análisis personal que interprete los resultados (2–3 líneas/ 2pts)

Nº	Pregunta / Requerimiento	Tipo de gráfico sugerido
1	Cantidad total de titulados por año	Gráfico de líneas
2	Top 10 regiones con más titulados (acumulado)	Barras horizontales
3	Evolución anual de titulados para Top 5 regiones	Líneas
4	Distribución de titulados por género	Torta o barras apiladas
5	Comparación mujeres vs hombres por año	Líneas
6	Top 15 instituciones con más titulados	Barras horizontales
7	Top 15 sedes con más titulados	Barras horizontales
8	Titulados por área del conocimiento	Barras horizontales
9	Top 20 carreras con más titulados	Barras horizontales
10	Clasificación de instituciones por niveles 1, 2 y 3	Barras agrupadas

Nota: Cada gráfico debe incluir un título claro, ejes etiquetados y fuente:

Fuente: Elaboración propia, datos de titulados 2007–2020.

## Parte III – Informe de análisis

Esta parte se evalúa en conjunto con la anterior

Deberá **documentar y explicar su proceso analítico** utilizando el enfoque **KDD (Knowledge Discovery in Databases)**, integrando los resultados obtenidos en el código.

Estructura mínima del documento PDF:

### 1. Portada:

- Nombre del estudiante
- Asignatura y evaluación
- Fecha
- Título: "*Análisis KDD – Titulados Educación Superior 2007–2020*"

### 2. Introducción:

Breve descripción del propósito del análisis y del conjunto de datos utilizado.

### 3. Desarrollo del proceso KDD:

Describa brevemente las etapas aplicadas:

- **Selección de datos:** Identificación del archivo y sus variables principales.
- **Limpieza de datos:** Transformaciones realizadas y correcciones aplicadas.
- **Transformación:** Creación de nuevas variables, agrupaciones o conversiones.
- **Minería de datos:** Resumen de los cálculos o consultas generadas (las 10 preguntas).
- **Interpretación y evaluación:** Principales hallazgos y observaciones relevantes.

### 4. Evidencias gráficas:

Incluya los **gráficos generados en el código R** con su respectivo número de pregunta, título, breve descripción e interpretación personal.

Ejemplo:

**Figura 1:** Titulados por año (2007–2020).

Se observa una tendencia creciente hasta 2018, con una leve baja en 2020.

### 5. Conclusiones:

- Principales descubrimientos y patrones observados.
- Reflexión sobre la utilidad del análisis de datos en la educación superior.
- Posibles líneas de análisis futuro (ejemplo: por provincia o edad promedio).