



Bases de Datos Masivas (11088)
Departamento de Ciencias Básicas

TRABAJO PRÁCTICO 02: Preprocesamiento de datos

-Análisis, Limpieza, Transformación e Integración-

Introducción:

En este trabajo se abordan cuestiones relacionadas con las tareas de Preprocesamiento de datos previo a la etapa de descubrimiento del conocimiento. Entre las tareas que se abordan, se encuentra la integración, limpieza, selección y transformación de variables, así como técnicas de reducción de dimensionalidad de un dataset, a efectos de reconocer aquellos atributos y escalas que mejor lo representan.

Se plantean ejercicios y datasets cuyas resoluciones serán realizadas mediante el lenguaje R.

Limpieza de datos:

1. *Datos faltantes.* Se cuenta con el dataset *encuesta_universitaria.csv*, el cual posee valores faltantes para la variable *tiempo_traslado*. Aplique los siguientes métodos a efectos de reemplazar esos valores:
 - a. Verifique cual es la proporción de valores faltantes respecto a la cantidad total de instancias del dataset.
 - b. Genere un nuevo atributo utilizando solo los registros con valores observados para el atributo.
 - c. Genere un nuevo atributo en el que sustituya los valores faltantes por la media encontrada para el atributo.
 - d. Genere un nuevo atributo en el que sustituya los valores faltantes de acuerdo al método de "hot deck imputation".
 - e. Analice los resultados encontrados a partir de la aplicación de los métodos anteriores. Compare los mismos realizando gráficos sobre los valores resultantes en cada caso.



Bases de Datos Masivas (11088)
Departamento de Ciencias Básicas

2. *Manejo de Ruido*. Para el dataset anterior, avance sobre las siguientes operaciones para los atributos numéricos (cuantitativos continuos):
 - a. Verifique en primer lugar la distribución de los datos, utilice algún método gráfico para esto.
 - b. Realice un suavizado utilizando *binning* por *frecuencias iguales* (con 5 bins) y estime el valor del bin por el cálculo de medias.
 - c. Ahora, realice el suavizado por *anchos iguales* (con 5 bins) y compare los resultados gráficamente.
3. *Detección de outliers*. Ahora, trabaje sobre el mismo atributo del dataset original con las siguientes consignas:
 - a. Verifique la existencia de *outliers* en el atributo *tiempo_traslado* en función del resto de los atributos. ¿En todos los casos se trata de un valor anómalo?
 - b. Aplique las técnicas de análisis y detección vistas en clase: IRQ, SD (seleccione el N que mejor se adapte a su criterio) y Z-Score (seleccione el umbral que mejor se adapte a su criterio).
 - c. Concluya respecto a los resultados obtenidos con cada técnica.

Reducción de dimensionalidad:

4. A partir del dataset *auto-mpg.data-original.txt*¹, se solicita trabajar sobre las siguientes consignas:
 - a. Evalúe la relación entre atributos a partir del coeficiente de correlación de Pearson y un análisis gráfico de heatmap² para estudiar la posibilidad de eliminar redundancia en el dataset. En caso de corresponder, aplique las técnicas de Reducing Highly Correlated Columns trabajadas en clase.
 - b. Verifique a través del Test de Chi-Cuadrado si existe dependencia entre pares de atributos discretos. Determine en qué casos es conveniente reducir dimensionalidad.

¹ Disponible en: <https://archive.ics.uci.edu/ml/datasets/Auto+MPG>

² Explore la instrucción *heatmap.2* de la librería *gplots*.



Bases de Datos Masivas (11088)
Departamento de Ciencias Básicas

5. *Análisis de Componentes Principales.* Cargue en R el dataset *europa.dat* y conteste las siguientes consignas a través de las funcionalidades provistas por esa herramienta:
- Calcule la matriz de covarianzas. ¿Qué nos indica la misma sobre los atributos del dataset?
 - Realice ahora el análisis de componentes principales. ¿Cuánto explica de la variación total del dataset la primera componente? ¿Y si se incorpora la segunda? ¿Y el primer auto-valor?
 - Grafique el perfil de variación de las componentes en un gráfico de dispersión donde las X es la componente y la Y la varianza.
 - Analice la matriz de loading. ¿Qué información provee? ¿Qué variables están más correlacionadas con la primera componente?
 - Genere un gráfico de biplot y explique brevemente que información le provee el mismo.
 - En función de los análisis realizados en los puntos anteriores. ¿Cuántas componentes principales elegiría para explicar el comportamiento del dataset? Justifique esa cantidad.

Transformación de datos:

6. *Discretización.* A partir del dataset *encuesta_universitaria.csv*, opere sobre el atributo *tiempo_traslado* de la siguiente manera:
- Transforme el atributo a discreto, definiendo 5 rangos de acuerdo al análisis de frecuencia de los valores encontrados para el atributo.
 - Transforme el atributo a discreto, definiendo 5 rangos de acuerdo al método de anchos iguales.
 - Transforme el atributo a discreto, definiendo usted, según su criterio, 5 rangos distintos con sus respectivas etiquetas.
 - Analice los resultados encontrados. Compare los mismos realizando gráficos de frecuencia sobre los intervalos resultantes en cada caso. ¿Qué conclusiones se pueden obtener en términos del balanceo de las mismas de acuerdo a la técnica utilizada?



Bases de Datos Masivas (11088)
Departamento de Ciencias Básicas

7. *Normalización*. Trabaje sobre las siguientes consignas:

- a. A partir del dataset *encuesta_universitaria.csv*, opere sobre el atributo *tiempo_traslado* de la siguiente manera:
 - i. Normalice el atributo utilizando la técnica de mínimo-máximo.
 - ii. Ahora, normalice el atributo mediante la técnica de z-score propuesta en el libro “Data Mining. Concepts & Techniques de Jiawei Han & otros”.
 - iii. Por último, utilice la técnica de escalado decimal para llevar adelante la tarea de normalización.
- b. Analice los resultados encontrados. Compare los mismos realizando gráficos sobre los atributos resultantes en cada caso.

Referencias sugeridas:

Principal component analysis. Hervé Adbi & otros. 2010.

Data mining and the impact of missing data. Marvin L. Brown & otros. 2003.

Data Mining. Concepts & Techniques. Jiawei Han and Micheline Kamber. 2006.