

A thick dark blue vertical bar runs along the left edge of the page. A purple arrow-shaped banner points to the right from this bar, containing the date '5-1-2025'. In the bottom-left corner, several thin, curved lines in dark blue and light grey sweep upwards and to the right.

5-1-2025

Captura y Preparación de Datos.

PEC4

Juan Luis Acebal Rico

GRADO DE CIENCIA DE DATOS APLICADA

Indice

Parte teórica	2
Q1.....	2
¿Qué diferencia hay entre selección o extracción de características y reducción de la dimensionalidad?.....	2
Q2.....	2
En los métodos de selección de características, ¿En qué se diferencian los métodos de envoltura de los métodos embebidos?.....	2
¿Y los filtros de los métodos embebidos?	2
Q3.....	2
¿Por qué se recomienda ser más conservadores con los filtros en el proceso de extracción de características?	3
¿Cómo funcionan los filtros?	3
¿Sirven los mismos filtros para un problema de regresión y de clasificación?	3
Q4.....	3
¿Cuáles son los métodos más conocidos para la codificación de las variables categóricas? Descríbelos brevemente y pon un ejemplo de cada uno.	3
Bibliografía	4

Parte teórica

Q1.

¿QUÉ DIFERENCIA HAY ENTRE SELECCIÓN O EXTRACCIÓN DE CARACTERÍSTICAS Y REDUCCIÓN DE LA DIMENSIONALIDAD?

Selección de características: se elige un subconjunto de las variables originales, sin modificar su significado (de ellas, aunque puede que del dataset ya que tiene menos columnas), por ejemplo si tenemos una variable que es booleana que representa el sexo pero luego tenemos una variable que representa sexo pero también LGTBI, podríamos quedarnos con la segunda (o la primera, depende el contexto) ya que tiene información repetida.

Extracción de características: se generan nuevas variables (transformaciones o combinaciones de las originales), frecuentemente con el objetivo de reducir la dimensionalidad, por ejemplo: patrimonio (que sería la columna activos - pasivos).

Reducción de la dimensionalidad: bajar el número de variables (dimensiones) mediante transformaciones que capturan la mayor parte de la información original (PCA es la más conocida aunque también hemos visto SVD, LDA,...).

Q2.

EN LOS MÉTODOS DE SELECCIÓN DE CARACTERÍSTICAS, ¿EN QUÉ SE DIFERENCIAN LOS MÉTODOS DE ENVOLTURA DE LOS MÉTODOS EMBEBIDOS?

Métodos de envoltura (wrapper): Entrenan un modelo buscando muchas combinaciones de características.

Métodos embebidos (embedded): La selección de características sucede dentro del proceso de entrenamiento del modelo.

¿Y LOS FILTROS DE LOS MÉTODOS EMBEBIDOS?

Filtros: aplican criterios estadísticos (correlación, chi-cuadrado, etc) sin depender de un modelo concreto.

Embebidos: la importancia de las características se determina dentro del propio entrenamiento del modelo.

Q3.

¿POR QUÉ SE RECOMIENDA SER MÁS CONSERVADORES CON LOS FILTROS EN EL PROCESO DE EXTRACCIÓN DE CARACTERÍSTICAS?

Los filtros descartan variables basándose en métricas más o menos simples, que además, por ejemplo, tengan poca correlación con una variable objetivo, pero luego puede darse el caso que en combinación con otras, SI que tengan relación con la variable objetivo.

¿CÓMO FUNCIONAN LOS FILTROS?

Se calcula una puntuación para cada característica que luego se puede seleccionar por un ranking o por un umbral

¿SIRVEN LOS MISMOS FILTROS PARA UN PROBLEMA DE REGRESIÓN Y DE CLASIFICACIÓN?

Si, pero se usan distintos puntajes, por ejemplo, Pearson, y chi cuadrado para regresión y clasificación respectivamente.

Q4.

¿CUÁLES SON LOS MÉTODOS MÁS CONOCIDOS PARA LA CODIFICACIÓN DE LAS VARIABLES CATEGÓRICAS? DESCRÍBELOS BREVEMENTE Y PON UN EJEMPLO DE CADA UNO.

One-Hot Encoding (*)

Crea una columna binaria para cada categoría

Ejemplo: "Color" con categorías {Rojo, Azul} se convierte en Color_Rojo, Color_Azul. El problema es que puede crear muchas columnas si hay muchas categorías., ya que si hay 10 colores para un coche, tendríamos 9 de ellos como false y 1 de ellos como true.

Label Encoding

Asigna un número entero a cada categoría.

Ejemplo: "Color" = {Rojo=0, Azul=1, Verde=2}.

Ordinal Encoding

Parecido a Label Encoding, pero para categorías con orden natural (p.ej., Tallas S < M < L).

Ejemplo: Talla = {S=1, M=2, L=3}. También como la parte práctica, education de 1 al 16 según el nivel de estudios....

Target (Mean) Encoding

Sustituye cada categoría por el promedio de la variable objetivo para esa categoría

Ejemplo: si Color=Rojo suele asociarse a un valor medio de 0.8 en la variable objetivo, se reemplaza por 0.8. Si education doctorado es en promedio un income de >50mil, que sería un 1, podríamos decir que education_doctorado sería cercano a 1. Siendo 1 en income >50mil y el número en education_doctorado sería el promedio de las personas que tienen doctorado en el valor de la variable objetivo.

Bibliografía

(*) <https://codificandobits.com/blog/como-codificar-datos-categoricos/>
<https://www.hackersrealm.net/post/target-or-mean-encoding-python>