

# 7506-2020-2 Examen por promocion

Santiago Locatelli

TOTAL POINTS

**62 / 100**

## QUESTION 1

### 1 Clustering 5 / 20

- ✓ - 5 pts b) mal.
- ✓ - 10 pts c) mal.

## QUESTION 2

### 2 Recomendaciones 20 / 20

- ✓ - 0 pts Correct

## QUESTION 3

### 3 PageRank 8 / 20

- ✓ - 12 pts Plantea pero no resuelve

## QUESTION 4

### 4 Streaming 20 / 20

- ✓ - 0 pts Correct

## QUESTION 5

### 5 Feature Engineering 9 / 20

Punto a

- ✓ - 2 pts Elimina datos importantes, como la fecha de publicación o la ciudad.
- ✓ - 3 pts Usa los features que vienen pero no considera agregar nuevos features.
- ✓ - 2 pts No procesa el campo Servicios.
- ✓ - 1 pts No aprovecha que el estado de la propiedad tiene un orden implícito para codificarlo con label encoding.
- ✓ - 2 pts Utiliza binary encoding para la ciudad sin evaluar codificarlo con mean encoding.

Punto b

- ✓ - 1 pts Falta detalle de qué campos debería normalizar y cómo.

1 Es la fecha de publicación

2

Es la fecha de publicación. Ud cree que el precio es el mismo si la propiedad se vendió el año pasado o hace 20 años?

3 Por?

4 Por?

5 Por qué binary encoding y no mean encoding?

6 Por qué?

7 El estado quizás sea mejor encodearlo con label encoding ya que hay un orden implícito en los datos.

- a) Sabiendo que el <sup>algoritmo de</sup> Lloyd agarra  $3(k)$  Puntos del Set de datos al azar y que hay 6 Puntos en el Set de datos, podemos calcular las configuraciones posibles. La primera vez que se elija pueden ser 6 opciones. La segunda vez que se elija pueden ser 5 opciones. " Tercera " " " " 4 opciones

Por lo que por regla de la cadena tenemos que las configuraciones posibles son:  $6 \cdot 5 \cdot 4 = \boxed{120}$

- b) Solo 1. Analizando los puntos dados, vemos que hay 3 Pares de Puntos, los cuales estos pares tienen distancia mínima entre ellos. Estos Pares son:

Par 1:  $x_1$  y  $x_4$  Par 2:  $x_2$  y  $x_5$  Par 3:  $x_3$  y  $x_6$

El hecho de que tome de a Pares se debe a la ~~proporción~~ Proporción entre la cantidad de Puntos y la cantidad Clusters.

Entonces, tenemos como conclusión que los clusters van a converger al promedio de cada Par.

Calculamos los Puntos.

$$\bullet (0, 3)$$

$$\bullet \left( \frac{8+8}{2}, \frac{6+0}{2} \right) = (8, 3)$$

$$\bullet \left( \frac{16+16}{2}, \frac{0+6}{2} \right) = (16, 3)$$

## 1 Clustering 5 / 20

✓ - 5 pts b) mal.

✓ - 10 pts c) mal.

Santiago Locatelli 104107.

2. Estimar la Probabilidad de que al usuario 6 le guste la Serie 'D'.

Lo Primero que hacemos es buscar los 3 Vecinos más cercanos al usuario 6.

Calculamos las semejanzas de Jaccard:

$$J(6,1) = \frac{0}{6} = 0 \quad J(6,2) = \frac{2}{6} = \frac{1}{3}$$

$$J(6,3) = \frac{2}{6} = \frac{1}{3} \quad J(6,4) = \frac{2}{6} = \frac{1}{3}$$

$$J(6,5) = \frac{1}{5}$$

Por lo tanto los ~~que~~ 3 Usuarios más semejantes a 6 son los Usuarios: 2, 3, 4

Como solamente se trata de Upvotes y Down votes, la Probabilidad de que le guste la serie está definida como el porcentaje de Upvotes sobre el total, es decir, que en este caso la Probabilidad de que le guste la Serie 'D' al usuario 6 es  $\frac{1}{3}$ .

Observación: si se usara la formula de ponderación, para obtener la Probabilidad, usando como 'Si' = 1 y 'No' = 0 da el mismo resultado.

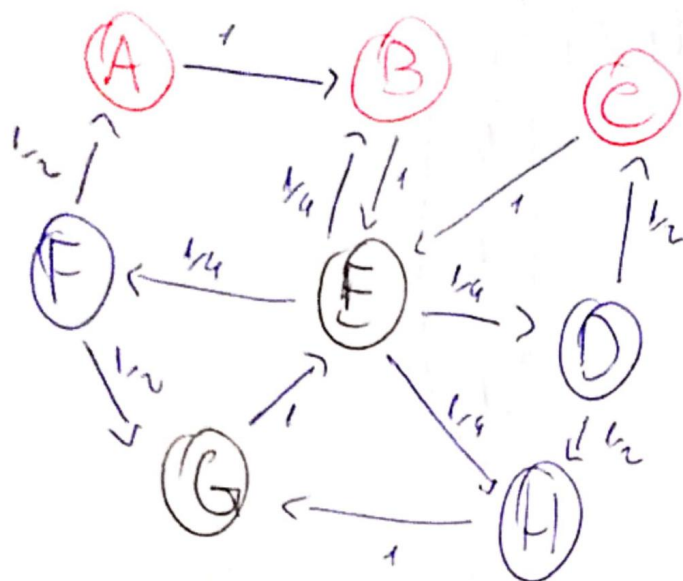
2 Recomendaciones 20 / 20

✓ - 0 pts Correct



3. ~~A continuación se muestra un grafo de transición.~~  
Teniendo el grafo:

Asumo un  $\beta = 0.8$ .



Podemos armar la matriz: A

	A	B	C	D	E	F	G	H
A	0	0	0	0	0	1/2	0	0
B	1	0	0	0	1/4	0	0	0
C	0	0	0	1/2	0	0	0	0
D	0	0	0	0	1/4	0	0	0
E	0	1	1	0	0	0	1	0
F	0	0	0	0	1/4	0	0	0
G	0	0	0	0	0	1/2	0	1
H	0	0	0	1/2	1/4	0	0	0

~~Y para darle importancia a la semántica de economía~~  
~~de la matriz la columna~~  
~~sumamos a~~

Entonces para calcular el Page Rank se debe calcular:

$$\text{PageRank} = \beta \cdot A \cdot \begin{pmatrix} 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \\ 1/8 \end{pmatrix} + (1-\beta) \cdot \begin{pmatrix} 0 \\ 0 \\ 0 \\ 1/3 \\ 0 \\ 1/3 \\ 0 \\ 1/3 \end{pmatrix}$$

No Resuelvo la cuenta por falta de tiempo.

### 3 PageRank 8 / 20

✓ - 12 pts Plantea pero no resuelve



9. a) Como por cada dato del Stream se le aplica la función de hashing y se suma 1 en el número de bucket que dio como resultado, la suma de todos los valores en los filtros debe ser la misma.

Notando que  $F_1 = F_2 = 20$ , tenemos que pedir que A o B sea igual a 20.

$A = 20$  y  $B = 21$ , por lo tanto el filtro correcto es el A.

b) Los Filtros nos quedan:

$$F_1 = [4, 0, 8, 0, 0, 0, 5, 3]$$

$$F_2 = [0, 0, 3, 3, 4, 4, 6, 0]$$

$$F_3 = [5, 0, 1, 3, 2, 4, 4, 1]$$

0 1 2 3 4 5 6 7

Para lograr la mayor estimación posible, los resultados de las funciones de hash deberían caer sobre los buckets que almacenan los valores más grandes.

Para lograr lo mencionado los valores de las funciones de hash deberían ser:

$$h_1 = 2$$

$$h_2 = 6$$

$$h_3 = 0$$

Entonces ~~la~~ la estimación es  $\min(8, 6, 5) = 5$ .

4 Streaming 20 / 20

✓ - 0 pts Correct

a) Asumo que el feature "fecha" es la fecha actual. 1

- Crearía una nueva feature que represente el año en <sup>que</sup> la propiedad se terminó (se estrenó), a través de la diferencia de la fecha con la antigüedad.

- Crearía también una feature que represente una relación entre los m<sup>2</sup> totales y los m<sup>2</sup> cubiertos. De la forma:

$$\frac{\text{m}^2 \text{ cubiertos}}{\text{m}^2 \text{ totales}}$$

Notar que sería un número menor o igual a 1. Y no precaria que m<sup>2</sup> totales sea cero, porque sino la propiedad no tendría sentido.

- Haría la concatenación de las features categóricas: Provincia y Ciudad, ya que están ultramente relacionadas.

Las features que no tendría en cuenta son:

- Población-ciudad. 3

- Cantidad-escuelas-cercanas.

- m<sup>2</sup>-totales
- m<sup>2</sup>-cubiertos
- Servicios.

} Es decir, las descartaría después de la creación de la feature ya mencionada. 4

- Fecha 2

También, es muy importante la etapa de encoding de las features que son categóricas.

Para nueva feature Provincia-ciudad utilizarla Binary Encoding, ya que es muy probable que haya muchos valores diferentes. 5

Para las features tipo-propiedad y estado-propiedad aplicaría One Hot Encoding. ~~Servicios y Baños~~ 6 7

## Ejemplo:

Año-estreno: 2010  
Relación\_mz\_cub\_total: 0.5.  
Habitaciones: 8  
Dormitorios: 4  
Baños: 2  
Cocheras: 1  
Provincia\_ciudad\_1: 0  
Provincia\_ciudad\_2: 0  
Provincia\_ciudad\_3: 1  
:  
Provincia\_ciudad\_La cantidad\_valores: 0  
Estado\_malo: 0  
Estado\_regular: 0  
Estado\_bueno: 1  
Antigüedad: 11

- b) Como en modelos basados en árboles las Features numéricas no afectan, pero en las redes ~~neuronales~~ neuronales si afectan drásticamente, realizó una normalización de las features numéricas. Más que nada a las features Año-estreno y antigüedad, que son las que más discrepan entre valores puede haber.

## 5 Feature Engineering 9 / 20

Punto a

- ✓ - 2 pts Elimina datos importantes, como la fecha de publicación o la ciudad.
- ✓ - 3 pts Usa los features que vienen pero no considera agregar nuevos features.
- ✓ - 2 pts No procesa el campo Servicios.
- ✓ - 1 pts No aprovecha que el estado de la propiedad tiene un orden implícito para codificarlo con label encoding.
- ✓ - 2 pts Utiliza binary encoding para la ciudad sin evaluar codificarlo con mean encoding.

Punto b

- ✓ - 1 pts Falta detalle de qué campos debería normalizar y cómo.

- 1 Es la fecha de publicación
- 2 Es la fecha de publicación. Ud cree que el precio es el mismo si la propiedad se vendió el año pasado o hace 20 años?
- 3 Por?
- 4 Por?
- 5 Por qué binary encoding y no mean encoding?
- 6 Por qué?
- 7 El estado quizás sea mejor encodearlo con label encoding ya que hay un orden implícito en los datos.