

1) Al leer el enunciado entiendo que el resultado no es satisfactorio ya que hay puntos débiles que son ni muy poderosos, ni muy débiles por lo tanto lo que se puede cambiar es la clasificación de dichos ^{as} débiles puntos. Para esto se me ocurren diversas soluciones, una podría ser aumentar la cantidad de clusters y en lugar de solo tomar dos clusters tomar 3 o 4 para que haya uno o dos clusters para estos casos límite.

Otra solución podría ser aplicar soft K-means el cual es una generalización de K-means que permite a cada punto pertenecer al mismo tiempo a más de un cluster. Se busca estimar la probabilidad de que el punto pertenezca a cada cluster. Con esta solución, los países que fueran más poderosos serían de un rango más claro y los más débiles de un rango claro. De esta forma podemos ver el grado de pertenencia de los puntos a cada cluster, y los que están en el medio no serían asignados ni a uno ni a otro necesariamente.

Otra posible alternativa sería utilizar DBSCAN, ya que este algoritmo permite tener puntos que no son parte de ningún cluster, por lo tanto si solo queremos estimar los poderosos y los débiles como posibles grupos los puntos que no pertenezcan a estas regiones limítrofes pueden ser descartados.

De las soluciones propuestas, la primera la descarto ya que entiendo que se busca estimar solo los grupos: fuertes y débiles.

* Propongo utilizar soft K means:

Ahora en lugar de original a cada punto me centrado más cercano como hace K means, se le asigna a cada punto dos probabilidades, dando cada una de ellas la probabilidad que el punto pertenezca a uno de los dos clusters, de tal manera que cuanto más cerca está el punto del centrado mayor es su probabilidad.

Para esto calculo la distancia del punto a cada centrado y la probabilidad de pertenecer a ~~este~~ un ~~cluster~~ cluster será: $1 - \frac{\text{dist al cluster}}{\text{suma de todas distancias}}$

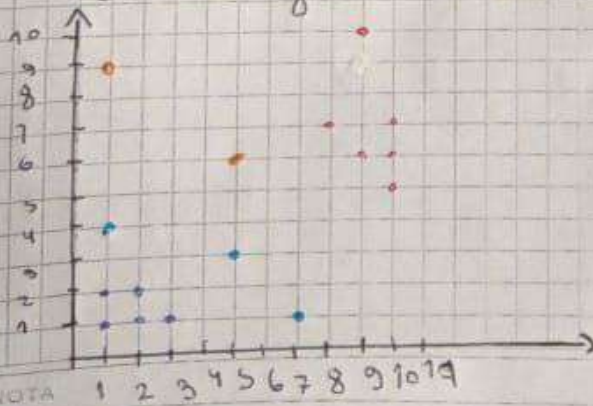
Utilizando por ejemplo la distancia euclídea como métrica. El algoritmo es el siguiente:

Calcular los centrados, originales lo voy como K means

Calcular la distancia de cada punto a los centrados

Calcular la probabilidad de pertenecer a c / cluster

La clustereación iterativa con este método será:



$$P(\text{Rosa}) = 1 - \frac{\text{dist Rosa}}{\text{dist Rosa} + \text{dist Angel}}$$

$$P(\text{Angel}) = 1 - \frac{\text{dist Angel}}{\text{dist Rosa} + \text{dist Angel}}$$

3
Considero que este resultado es más apropiado ya que los países que antes eran considerados como poderosos o débiles, pero en realidad no son ninguno de estos, ahora son ~~considerados~~ ^{considerados} por su grado de pertenencia a cada cluster, obteniendo una clusterización más interesante y una resolución más rica.

Por ejemplo el punto (2, 3) antes era designado como poderoso, aunque hubiera mandado solo 1 delegación, los países que mandaban pocos atletas no suelen ser muy poderosos, pero aun así ganaron 9 medallas, por lo que lo vemos mejor, es un país más poderoso que otros pero no es de los más poderosos.

Por esta razón considero que la solución propuesta es mejor que la de DB con ya que en lugar de eliminar los datos en el medio de los dos grupos, los mantenemos y obtenemos información interesante sobre ellos.

Donde Raimundo
102848

Promocional

4

2) Utilizar CF User-User con $K=2$ para estimar la calificación del item 3 para el usuario a.

En primer lugar, buscar los dos usuarios más similares al usuario a.

$$a: [1 \ ? \ ? \ 5 \ 1 \ ? \ 2 \ 3] \rightarrow \bar{a} = 2,4$$

$$b: [3 \ 2 \ ? \ 1 \ 1 \ ? \ 3 \ 2] \rightarrow \bar{b} = 2$$

$$c: [2 \ ? \ 4 \ ? \ 2 \ ? \ 4 \ 2] \rightarrow \bar{c} = 2,8$$

$$d: [1 \ 2 \ ? \ ? \ 2 \ ? \ ? \ 3] \rightarrow \bar{d} = 1,75$$

$$e: [? \ 3 \ 3 \ 1 \ 1 \ 4 \ ? \ 2] \rightarrow \bar{e} = 2,6$$

$$f: [4 \ ? \ 5 \ 1 \ ? \ ? \ 4 \ 1] \rightarrow \bar{f} = 3$$

Los datos faltantes no los voy a tener en cuenta a la hora de calcular los valores

Calcular el promedio de cada usuario y restar esto al vector para centrar las calificaciones

$$a: [-1,4 \ ? \ ? \ 2,6 \ -1,4 \ ? \ -0,4 \ 0,6]$$

$$b: [1 \ 0 \ ? \ -1 \ -1 \ ? \ 1 \ 0]$$

$$c: [-0,8 \ ? \ 1,2 \ ? \ -0,8 \ ? \ 1,2 \ -0,8]$$

$$d: [-0,75 \ 0,25 \ ? \ ? \ -0,75 \ ? \ ? \ 1,25]$$

$$e: [? \ 0,4 \ 0,4 \ ? \ -1,6 \ 1,4 \ ? \ -0,6]$$

$$f: [1 \ ? \ 2 \ -2 \ ? \ ? \ 1 \ -2]$$

Calcular similitud con correlación de Pearson

$$\text{sim}(x, y) = \frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)^2} \sqrt{\sum_{s \in S_{xy}} (r_{ys} - \bar{r}_y)^2}}$$

* Similitud entre a y b: solo tengo en cuenta puntos en común

$$A: (-1,4 \ 2,6 \ -1,4 \ -0,4 \ 0,6) \quad \|A\| = 3,35$$

$$B: (1 \ -1 \ 1 \ 1 \ 0) \quad \|B\| = 2$$

Normalizando:

$$A: (-0,418, 0,776, -0,418, -0,119, 0,179)$$

$$B: (0,5, -0,5, 0,5, 0,5, 0)$$

$$\text{Sim}(A, B) = A \cdot B = -0,8655$$

* Similitud entre a y c:

$$A: (-1,4, -1,4, -0,4, 0,6) \quad \|A\| = 2,11$$

$$C: (-0,8, -0,8, 1,2, -0,8) \quad \|C\| = 1,83$$

Normalizando:

$$A: (-0,66, -0,66, -0,19, 0,284)$$

$$C: (-0,437, -0,437, 0,656, -0,437)$$

$$\text{Sim}(A, C) = A \cdot C = 0,328$$

* Similitud entre a y d:

$$A: (-1,4 \ -1,4 \ 0,6 \ 1) \quad \|A\| = 2,02$$

$$D: (-0,75 \ -0,75 \ 1,25) \quad \|D\| = 1,64$$

Normalizando:

$$A: (-0,676, -0,676, 0,29)$$

$$D: (-0,457, -0,457, 0,762)$$

$$\text{Sim}(A, D) = 0,839$$

* similitud entre α y B :

$$A: (-1.4 \quad 0.6) \quad \|A\| = 1.52$$

$$E: (-1.6 \quad -0.6) \quad \|E\| = 1.71$$

Normalizado:

$$A: (-0.92, 0.395)$$

$$E: (-0.936, -0.35)$$

$$\text{sim}(A, E) = 0.723$$

* similitud entre α y F :

$$A: (-1.4 \quad 2.6 \quad -0.4 \quad 0.6) \quad \|A\| = 3.04$$

$$F: (1 \quad -2 \quad 1 \quad -2) \quad \|F\| = 3.16$$

Normalizado:

$$A: (-0.46, 0.855, -0.132, 0.197)$$

$$F: (0.316, -0.633, 0.316, -0.633)$$

$$\text{sim}(A, F) = -0.85$$

- Los dos usuarios más similares a A son D y E
- Sin poder calificar al usuario α por el ítem 3 se:

$$r_{u, i} = \frac{\sum_{j \in N} S_{x, j} \cdot r_{y, j}}{\sum_{j \in N} S_{x, j}}$$

$$r_{u, i_3} = \frac{0.723}{0.723 + 0.839} = \frac{0.839 \cdot 3}{0.839} = 3$$

- Como el usuario D no califica el ítem 3, no lo tengo en cuenta para el cálculo
- El usuario α califica con 3 estrellas el ítem 3.

3)

$$P_A(0) = 0,2 \cdot \frac{1}{8}$$

$$P_A(1) = 0,8 \left(\frac{1}{2} P_{A0} + P_{A2} \right) + 0,2 \cdot \frac{1}{8}$$

$$P_A(2) = 0,8 \left(\frac{1}{2} P_{A0} + \frac{1}{4} P_{A3} \right) + 0,2 \cdot \frac{1}{8}$$

$$P_A(3) = 0,8 (P_{A8}) + 0,2 \cdot \frac{1}{8}$$

$$P_A(4) = 0,8 \left(P_{A1} + \frac{1}{4} P_{A3} \right) + 0,2 \cdot \frac{1}{8}$$

$$P_A(5) = 0,8 \left(\frac{1}{2} P_{A4} + \frac{1}{4} P_{A3} \right) + 0,2 \cdot \frac{1}{8}$$

$$P_A(6) = 0,8 \left(\frac{1}{4} P_{A3} \right) + 0,2 \cdot \frac{1}{8}$$

$$P_A(8) = 0,8 \left(\frac{1}{2} P_{A4} + P_{A5} + P_{A6} \right) + 0,2 \cdot \frac{1}{8}$$

a) Con estas formulas para calcular el Page Rank con teletransportacion haremos en el guiso, como en algoritmo y obtiene la siguiente tras 10 iteraciones

	0	1	2	3	4	5	6	8
0	1/8	1/8	1/8	1/8	1/8	1/8	1/8	1/8
...								
10	0,025	0,035	0,03	0,22	0,14	0,13	0,02	0,24

Redondeo

→ Para que sume 1

b) sabiendo que la page 0 no es confiable, vamos a aplicar Trust Rank y buscaremos en la teletransportacion a los paginas confiables:

$$P_A(0) = 0 \rightarrow \text{No tiene links de entrada y no buscaremos teletransportacion}$$

$$P_A(1) = 0,8 \left(\frac{1}{2} P_{A0} + P_{A2} \right) + 0,2 \cdot \frac{1}{7} \rightarrow \text{Confianza de Paginas confiables}$$

$$P_A(2) = 0,8 \left(\frac{1}{2} P_{A0} + \frac{1}{4} P_{A3} \right) + 0,2 \cdot \frac{1}{7}$$

$$P_A(3) = 0,8 (P_{A8}) + 0,2 \cdot \frac{1}{7}$$

$$P_A(4) = 0,8 \left(P_{A1} + \frac{1}{4} P_{A3} \right) + 0,2 \cdot \frac{1}{7}$$

$$P_A(5) = 0,8 \left(\frac{1}{2} P_{A4} + \frac{1}{4} P_{A3} \right) + 0,2 \cdot \frac{1}{7}$$

$$P_A(6) = 0,8 \left(\frac{1}{4} P_{A3} \right) + 0,2 \cdot \frac{1}{7}$$

$$P_A(8) = 0,8 \left(\frac{1}{2} P_{A4} + P_{A5} + P_{A6} \right) + 0,2 \cdot \frac{1}{7}$$

NOTA

Dante Remundo
102848

Promoumbe

Hoy 1 8
P.O.N.

Con estas nuevas modificaciones, el pago con stream
tus diez iteraciones es:

	0	1	2	3	4	5	6	8
0	$1/8$	$1/8$	$1/8$	$1/8$	$1/8$	$1/8$	$1/8$	$1/8$
...								
10	0	0,09	0,07	0,23	0,15	0,13	0,07	0,26 \rightarrow Runa 1

4) Stream: $\{3, 6, 6, 3, 3, 6, 5, 6\}$

$$h(x) = x \bmod 32$$

Flajolet Martin: unico estimador

\hookrightarrow > bits, se consideran bits a derecha

Primero calculo momento de orden 0 del stream

$M^0(5) = 3 \rightarrow$ se tienen tres elementos diferentes en el
stream: 3, 6 y 5

$$h(3) = 3 \bmod 32 = 3 \rightarrow 00011$$

$$h(5) = 5 \bmod 32 = 5 \rightarrow 00101$$

$$h(6) = 6 \bmod 32 = 6 \rightarrow 00110 \rightarrow r = 1$$

Por lo tanto el momento de orden 0 para el stream
se calcula como $2^r = 2^1 = 2$.

El algoritmo de Flajolet Martin estima con esta funcion
que el momento de orden 0 para el stream dado, es decir, la
cantidad de elementos diferentes es 2

Considero que la aproximacion es muy cercana a 3,
por lo que esta funcion pareciera andar bien al menos con la
peor cantidad de bits que tenemos. A su vez utilizando
un unico estimador solo podremos obtener resultados que sean
potencias de 2 y jamas obtendremos 3 como resultado.

B) Considero que plantear una función de hashing alternativa me mejorará sustancialmente los obtenidos en el punto a ya que con un único estimador podemos obtener un único bit 0 a la derecha o a la suma dos (en caso de haber más la estimación será par) con estos resultados podremos obtener 2 o 4 como momentos de orden 0 del stream, y los resultados distorsionan en una unidad del momento de orden 0 real que es 3. Además, con un único estimador solo podemos aproximar con valores que sean potencias de 2. En lugar de una función alternativa propondría utilizar varias funciones del estilo $(ax+b) \bmod 32$ de manera de obtener varios estimadores y repartirlos en grupos. Luego obtengo la máxima de los momentos de cada grupo y le aplico el promedio. De esta manera podrá aproximar a momentos que no sean potencias de 2 e incluso pueden ser impares.

C) Plantear función del tipo $h(x) = (ax+b) \bmod 32$ que emplee parámetros

$$h(x) = (0 \cdot x + 16) \bmod 32$$

$$h(3) = 16 \rightarrow 10000$$

$$h(5) = 16 \rightarrow 10000$$

$$h(6) = 16 \rightarrow 10000 \rightarrow p=4$$

$$M_0(5) = 2^4 = 16$$

Donato Benavides
102848

Promocional

10

Esta función es muy útil ya que genera muchos
colores al darle valores constantes y además
genera muchos bits en 0 a la derecha, por lo que el
algoritmo de Floyd-Martin produce un valor muy
cerca del momento de orden 0 del stream.

Dante Remando
102848

Promocional

11

5) Proponer datasets de train y/o test de al menos 4 puntos cada uno, y un modelo tal que se cumpla

a) una regresión cuyo MSE en test sea 0 pero en train sea 0,5

Propongo una regresión lineal: $\hat{y} = m \cdot x + b$

$$MSE: (\hat{y}_{\text{test}} - y_{\text{test}})^2$$

Set de Train:

x	y	$\hat{y}(x)$
0	0	$\sqrt{2}/2$
0	0	$\sqrt{2}/2$
0	0	$\sqrt{2}/2$
0	0	$\sqrt{2}/2$

$$\frac{1}{N} \sum (\hat{y} - y)^2 = 0,5$$

$\hookrightarrow N=4$, 4 puntos

$$\sum (\hat{y} - y)^2 = 2$$

$$(\hat{y} - y)^2 = \frac{\sqrt{2}}{2}$$

Set de Test:

x	y	$\hat{y}(x)$
$\sqrt{2}/2$	$\sqrt{2}/2$	$\sqrt{2}/2$
$\sqrt{2}/2$	$\sqrt{2}/2$	$\sqrt{2}/2$
$\sqrt{2}/2$	$\sqrt{2}/2$	$\sqrt{2}/2$
$\sqrt{2}/2$	$\sqrt{2}/2$	$\sqrt{2}/2$

$$\frac{1}{N} \sum (y - \hat{y}(x))^2 = 0$$

De esta manera $\hat{y}(x) = 0 \cdot x + b$

$$\hat{y}(x) = \sqrt{2}/2$$

$$EMST_{\text{train}} = \frac{1}{4} \cdot 4 \left(0 - \frac{\sqrt{2}}{2} \right)^2 = 0,5$$

$$EMST_{\text{test}} = \frac{1}{4} \cdot 4 \left(\frac{\sqrt{2}}{2} - \frac{\sqrt{2}}{2} \right)^2 = 0$$

Dante Renuelo
102848

Promocional

12

C) Arbol de decision entrenado sobre dataset de train
cuyas variables son el linery encontrado de una categoria
de 3 clases y tiene 0,75 de accuracy entren

Accuracy: $\frac{\text{hechos correctos}}{\text{elementos clasificados}}$

Modelo: ID3

Clase A: 00

Clase B: 01

Clase C: 10

} Binary encoding

Clase D: 11

↳ Label 0

Dataset:

Bit 0	Bit 1	Label
0	0	1
0	1	0
1	0	1
1	0	1

Bit 0	Bit 1
1	0
1	100%
1	1
0	1

Con este modelo al tratar de predecir el target
de la clase D sabemos que es 1 ya que su bit 0 es 1,
por esto con los 4 casos, reemplazaremos
bien 3 de ellos y mal 1, obteniendo un accuracy del
75%.