# 4A)

| = | ≠ | CARA | DOC |
|---|---|------|-----|
| 0 | 5 | 0 | 9 |
| 4 | 6 | 5 | 10 |
| 2 | 8 | 11 | 33 |
| 0 | 9 | 19 | 38 |
| 4 | 2 | 28 | 55 |
| 2 | 7 | 30 | 69  5 |
| 0 | 10 | 37 | 76 |

LAS POSIC. EN LOS DOCS ESTÁN INTERNAS DESDE
POS 1, Y TODO CODIFICADO EN $\gamma$

MADRE IGUERA GISTRADO MACEFICIO MANANTIAL MANIFIESTO
0        5       11      19        30      37
                         82        28

LÉXICO CONCATENADO CON FRONT
→ APLICANDO CODING PARCIAL
USANDO $n=3$ (CADA 3 TÉRMINOS
ESCRIBE EL TÉRMINO COMPLETO)

1R1 00 100 111 10 0 100 01 10101010100 ...
0      10                 33    38              55           69    76

010101110101000100

10100100110110101001001010101011111010101111010101101010100010010100
0      10                         33      38                    55

010001110100010101101010100100
69              76

→ ESTOS SON LOS PUNTEROS A DOCS CODIFICADOS COMO GAMMA, PERO TAMBIÉN INCLUYE INFO SOBRE
POSICIONES DE CADA TÉRMINO Y FRECUENCIA ( DOC 1 FREC 1 P1 ... PN DOC2 ... )

→ CON LOS DATOS DE LA TABLA PASADA Y DEL LÉXICO CONCATENADO, PUEDO OBTENER LOS TÉRMINOS

MADRE → MADRIGUERA → MAGISTRADO → MACEFICIO → MALETA → MANANTIAL → MANIFIESTO

→ P/ OBTENER ESTOS TÉRMINOS, SE REALIZA UNA BÚSQUEDA BINARIA SOBRE
EL ÍNDICE Y SE VA ACCEDIENDO A CADA POSICIÓN DEL LÉXICO Y LEYENDO TANTOS CARACTERES
IGUALES O DISTINTOS DEL TÉRMINO ANTERIOR SEGÚN ME DIGA CADA POSICIÓN DEL ÍNDICE

AHORA TENGO QUE OBTENER LOS DOCUMENTOS, FRECUENCIA Y POSICIONES DE CADA TÉRMINO SEGÚN
LOS PUNTEROS CONCATENADOS → COMO ESTÁN CODIFICADOS EN GAMMA, SÉ QUE UNA FORMA DE DECO-
DIFICAR GAMMA ES → LEER LA CANTIDAD DE CEROS CON LA QUE EMPIEZA CADA DECIMAL DE NÚMERO
→ CONCATENACIÓN

$$\gamma(X) = LEN(BINARIO(X)) - 1 \text{ CEROS} + BINARIO(X)$$

→ TODAS LAS DISTANCIAS (ENTRE DOCS VA
RESIDUALMENTE CODIFICADAS EN GAMMA

ENTRE LAS POSICIONES 0 y 10 → 1 0 1 0 0 1 0 0 00 → DOC 1, FREC 2, POSICIÓN 2 y 3
                                         1  2  2   3

ENTRE LAS POSICIONES 10 y 33 → 0 1 1 0 1 0 1 00 100 10 1010 1 1111 → DOC 3, FREC 2, POS 1 y 4
                                          3    2 1  4   1 2 1 3 11        → DOC 4, FREC 2, POS 1 y 8
                                                                         → DOC 5, FREC 1, POS 1

ENTRE POS 33 y 38 → 01011 0 → DOC 2, FREC 1, POS 1
                        2  11

ENTRE POS 38 Y 55 - 111 0101 011 010 1010 → DOC 1, FREC 1, POS 1
  1 1 1   2   3   2 1  2            → DOC 3, FREC 1, POS 3
                                    → DOC 5, FREC 1, POS 2

→ ENTRE 55 Y 69 → 0100100 100 0100 → DOC 2, FREC 2, POS 2 y 04
        2    2   2   4

→ ENTRE 69 Y 76 → 0101 0111 010 → DOC 3, FREC 1, POS 2
                    3   1   2

→ POSICIONES MÁS ALLA DE 76 → 01010110 1010 0100 → DOC 2, FREC 1, POS 3
                              2 1  3   2 1  4      → DOC 4, FREC 1, POS 4

[NRO LOS DOCUMENTOS] → SON 5 DOCUMENTOS

 MALDICE ... EN: DOC 1, FREC 2, POS 2 y 3
MAYORÍQUEZA EN : DOC 3, FREC 2, POS 1 y 4 ; DOC 4, FREC 2, POS 1 y 3 ; DOC5, FREC 1, POS 1

MAGISTRADO EN: DOC 2, FREC 1, POS 1

MALEFICO EN : DOC 1, FREC 1, POS 1 ; DOC 3, FREC 1, POS 3 ; DOC5, FREC 1, POS 2

MALETA EN : DOC 2, FREC 2, POS 2 y 4

MANUAL EN : DOC 3, FREC 1, POS 2

MANIFIESTO E : DOC 2, FREC 1, POS 3 ; DOC 4, FREC 1, POS 4


 D1 : MALEFICO MALDICE MALDICE
 D2 : MAGISTRADO, MALETA, MANIFIESTO, MALETA
 D3 : MAYORÍQUEZA, MANUAL, MALEFICO, MAYORÍQUEZA
 D4 : MAYORÍQUEZA, MAYORÍQUEZA, MANIFIESTO          HABRÁ ALGÚN ERROR
 D5 : MAYORÍQUEZA, MALEFICO,                         YA QUE NO APARECE LA POS
                                                     2 DEL DOC 4.

TF-IDF → TF: TERM FREQUENCY (CANT DE VECES QUE APARECE CADA TERMINO DE Q
         EN CADA DOC
      IDF → BUSCA DARLE UN PESO INVERSAMENTE PROPORCIONAL A SU FREC
         (LOS TERMINOS QUE APARECE EN MUCHOS DOCS NO A SER TERMOS IMPORTANTES
         QUE LOS Q' APARECE EN POCOS DOCS)

Q = "MAYORÍQUEZA MALEFICO"          IDF = $\log\left(\frac{N+1}{FT_i}\right)$ → N: CANT TOTAL DE DOCS)
                                                             → CANT DE DOCS QUE POSEE
                                                                EL TERMINO

|           | TF |    |    |    |    | IDF              |
|-----------|----|----|----|----|----|------------------|
|           | D1 | D2 | D3 | D4 | D5 |                  |
| MAYORÍQUEZA | 2  | 0  | 2  | 2  | 1  | log (6/3) = 0,301 |
| MALEFICO   | 1  | 0  | 1  | 0  | 1  | log (6/3) = 0,301 |

PARA RANQUEAR, CALCULAMOS → $R(Q,D) = \sum FT_i \cdot IDF$  (PARA CADA TERMINO DE Q Y EN CADA DOC)
                                        4°
D1 → 0 * 0,301 + 1*0,301 = 0,301 ; D2 → 0 ; D3 → 2*0,301 + 0,301 = 0,903 ; D4 → 2*0,301 = 0,602
                                                                    1°                    2°
          3°
D5 → 1*0,301 + 1*0,301 = 0,602    RANKING → 1°: D3, 2°: D4, 3°: D5, 4°: D1, MALDICE

① Ⓐ   MUERTES = SC. TEXTFILE ('MUERTES .CSV')
    PORCENTAJE_RAZA = SC.TEXTFILE ('RAZAS. CSV')

# FILTRO LAS MUERTES POR RAZA NEGRA
MUERTES_NEGROS = MUERTES. FILTER (LAMBDA X: X[3] == "BLACK")\
. MAP (LAMBDA X: (X[4], 1))    → (ESTADO, 1)

# SUMO LAS MUERTES
MUERTES_NEGROS = MUERTES_NEGROS. REDUCEBYKEY (LAMBDA X,Y: X+Y)

# ME QUEDO CON EL ESTADO DE MAYOR CANTIDAD DE MUERTES YA QUE ESE VA A TENER EL
MAYOR PORCENTAJE DE MUERTES

MUERTES_NEGROS. REDUCE (LAMBDA X,Y: X IF X[1] > Y[1] ELSE Y)

~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~

Ⓑ #  (STATE, (RAZA NEGRA + 1))
           SURCE.BLACK
                                                    • SUMO LOS PORCENTAJE
                                                    DE CADA CLASE Y EL
RDD1 = PORCENTAJE_RAZA. MAP (LAMBDA X: (X[0], (X[3],1))    CONTADOR
RDD1 = RDD1. REDUCEBYKEY (LAMBDA X,Y: (X[0]+Y[0], X[1]+Y[1]))

# CALCULO EL PROCESO POR ESTADO          → PROCESO

RDD1 = RDD1. MAP (LAMBDA X: (X[0], X[1][0] / X[1][1]))

# CON ESTOS DATOS ESTE RDD YA CALCULO EL PORCENTAJE DE LAS MUERTES
MUERTES_NEGROS = MUERTES_NEGROS. MAP (LAMBDA X: (X[0], X[1]/MUERTES_NEGROS. COUNT()))
# JOINEO RDD1 CON MUERTES_NEGROS

RDD2 = RDD1.JOIN (MUERTES_NEGROS) → # (ESTADO, (PORCENTAJE_MUERTES, PROCESO))

# AHORA CALCULA LA DIFERENCIA ENTRE PORCENTAJE Y PROCESO

$RDD2 = RDD2, MAP(LAMBDA\ X: (X[0], ABS(X[1][0] - X[1][1])))$

ESPACIO      DIFERENCIA

## AHE JUEGO CON EL TOP10

$RDD2. TAKE ORDERED (10, LAMBDA\ X: -X[1])$

② A) IMPORT PANDAS AS PD

METRICAS = PD. READ_CSV ('METRICS. CSV')
CLIENTES = PD. READ_CSV (' MEARCAN 'CLIENTS. CSV')

# ME QUEDO CON CLIENTES CON EL NUMERO DE CUENTA MAYOR AL DADO

CLIENTES. FILTER = CLIENTES [CLIENTES ['ACCOUNT_NUMBER'] > 25679247]

# ME QUEDO CON LAS METRICAS PEDIDAS

METRICAS_PEDIDAS = METRICAS [ (METRICAS ['METRIC'] == 'AWS.VPC.NETWORK_IN') \
        OR (METRICAS ['METRIC'] == 'AWS.VPC.NETWORK_RATE') \
        OR (METRICAS ['METRIC'] == 'AWS.VPC.NETWORK_OUT') ]

# HAGO MERGE 'INNER' POR CLIENT_ID

MERGES = CLIENTES_FILTER. MERGE (METRICAS_PEDIDAS, ON = 'CLIENT_ID')

# AGRUPO POR METRICA y ACCOUNT_NUMBER

MERGES = MERGES. GROUPBY ([ 'METRIC', 'ACCOUNT_NUMBER' ]). AGG ({'VALUE': 'MEAN'}) RESET INDEX()
        VALUE INDEX (ACCOUNT NUMBER)

MERGES. UNSTACK (). T. RENAME (COLUMNS = { 'AWS.VPC.NETWORK_IN': 'AWS.VPC. NETWORK_IN_MEAN',
        'AWS.VPC.NETWORK_OUT': 'AWS.VPC.NETWORK_OUT_MEAN',
        'AWS.VPC.NETWORK_RATE': 'AWS.VPC.NETWORK_RATE_MEAN'})

③ B)

$$v_1 = [1,5,3,2], v_2 = [0,0,1,2], v_3 = [4,4,5,0], v_4 = [5,1,0,1]$$

CALCULAR TECNICA DE HIPERPLANOS con $b=1$, $r=4$. HALLAR HIPERPLANOS ADECUADO PARA QUE EL + SUMCN $g$ A $v_1$ SEA $v_3$

CALCULO DE MH → $r \times b = 4$

PARA QUE LOS 2 VECS SEAN SOLUCION, TIENE QUE DAR HALLAR EL MISMO VALOR DE PRODUCTO MISMO ENTRE LOS VECS Y LOS HIPERPLANOS

$$u_1 = [1,-1,-1,1] →$$

| | | |
|---|---|---|
| $v_1 \cdot u_1 < 0$ | → MH1 = 1 | |
| $v_2 \cdot u_1 > 0$ | → MH2 = -1 | } $v_1$ SUMA A $v_3$ |
| $v_3 \cdot u_1 < 0$ | → MH3 = 1 | |
| $v_4 \cdot u_1 > 0$ | → MH3 = -1 | |

$$u_2 = [-1,1,1,1] →$$

| | | |
|---|---|---|
| $v_1 \cdot u_2 > 0$ | MH1 = 1 | |
| $v_2 \cdot u_2 < 0$ | MH2 = -1 | } $v_1$ SUMA A $v_3$ |
| $v_3 \cdot u_2 > 0$ | MH3 = 1 | |
| $v_4 \cdot u_2 < 0$ | MH4 = -1 | |

$$u_4 = [1,1,1,-1] → v_1 \cdot u_4 > 0$$
$$v_2 u_4 < 0$$
$$v_3 u_4 > 0$$
$$v_4 u_4 > 0$$

$$u_3 = [-1,-1,-1,-1] →$$

| | |
|---|---|
| $v_1 \cdot u_3 < 0$ | |
| $v_2 \cdot u_3 < 0$ | } → CON u4 ENCUENTRO 2 HIPERPLANOS ACERTADOS |
| $v_3 \cdot u_3 < 0$ | QUE HACE QUE u1 y u3 SEAN SOLUCION, |
| $v_4 \cdot u_3 < 0$ | NO CREAN EL RESULTADO |

| | MH1 | MH2 | MH3 | MH4 |
|---|---|---|---|---|
| $v_1$ | 1 | | | |
| $v_2$ | | | | |
| $v_3$ | | | | |
| $v_4$ | | | | |

| | u1 | u2 | u3 | u4 |
|---|---|---|---|---|
| u1 | 1 | -1 | 1 | -1 |
| u2 | 1 | -1 | 1 | -1 |
| u3 | -1 | -1 | -1 | -1 |
| u4 | 1 | -1 | 1 | 1 |

$v_1$ y $v_3$ SOLUCION, EL RESTO SUS MH USANDO

$$u_1 = [1,-1,-1,1]$$
$$u_2 = [-1,1,1,-1]$$
$$u_3 = [-1,-1,-1,-1]$$
$$u_4 = [1,1,1,-1]$$