# EE412 Foundation of Big Data Analytics, Fall 2021
# HW3

Name: Cao Viet Hai Nam

Student ID: 20200817

Discussion Group (People with whom you discussed ideas used in your answers):


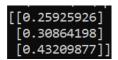On-line or hardcopy documents used as part of your answers:


## Answer to Problem 1:
## Exercise 5.1.2.

Code:

```
1.     import numpy as np
2.     beta = 0.8
3.     M = np.array([[1/3,1/2,0],
4.                   [1/3,0,1/2],
5.                   [1/3,1/2,1/2]])
6.     v = np.array([[1/3], [1/3], [1/3]])
7.     teleport_fac = (1-beta)*np.ones((3, 1)) / 3
8.     n = 50
9.     for i in range(n):
10.            v = beta*M@v + teleport_fac
11.    print(v)
```

**RESULT:**

```
[[0.25925926]
 [0.30864198]
 [0.43209877]]
```

Explanation:
        Line 2: initialize beta
        Line 3-5: Transition matrix
        Line 6: initial vector v will have 1/n for each component
        Line 7: Using taxation allow each suffer a small probability (1-beta) to teleport to a random page.
        Line 8: Number of iterations (about 50 will ensure for v to converge)
        Line 9-10: Iterate and update v.
        Line 11: print the result

# Exercise 5.3.1:

Code:

```
1.    import numpy as np
2.    beta = 0.8
3.    N = np.array([[0, 1/2, 1, 0 ],
4.                  [1/3, 0, 0, 1/2],
5.                  [1/3, 0, 0, 1/2],
6.                  [1/3,1/2,0, 0 ]])
7.    v1 = np.array([[1/4], [1/4], [1/4], [1/4]])
8.    v2 = np.array([[1/4], [1/4], [1/4], [1/4]])
9.    teleport_fac_A_only = (1-beta)*np.array([[1], [0], [0], [0]])
10.   teleport_fac_A_C = (1-beta)*np.array([[1/2], [0], [1/2], [0]])
11.   n = 50
12.   for i in range(n):
13.           v1 = beta*N@v1 + teleport_fac_A_only
14.           v2 = beta*N@v2 + teleport_fac_A_C
15.   print("A only: ")
16.   print(v1)
17.   print("A and C: ")
18.   print(v2)
```

**RESULT: (Page tank vector for both case)**

```
A only =
[[0.42857143]
 [0.19047619]
 [0.19047619]
 [0.19047619]]
A and C =
[[0.38571429]
 [0.17142857]
 [0.27142857]
 [0.17142857]]
```

Explanation:
  Line 2: initialize beta
  Line 3-6: Transition matrix
  Line 7-8 : initial vector v1, v2 will have 1/n for each component
  Line 9, 10: Using taxation allow each suffer a small probability (1-beta) to
       teleport to teleport set (only A or A and C)
  Line 11: Number of iterations (about 50 will ensure for both v to converge)
  Line 12-14: Iterate and update v1 and v2
  Line 15-18: print the result

**Answer to Problem 2**

② Mining Social - Network Graphs

<u>Exercise 10.3.2</u>

a) let $G$ be the bipartite graph.

let the "item" be the nodes on one side of $G$, we call it left side. Assume instance of $K_{s,t}$ has $t$ nodes on the left side. "baskets" corresponds to the nodes on the Right side.

let support threshold be $s$ = # nodes that $K_{s,t}$ has on the right side

⟹ frequent itemset of size $t$ and $s$ of the baskets in which all those items appear form $K_{s,t}$

Suppose the degree of $i^{th}$ node on the Right side is $d_i$ (size of $i^{th}$ basket)

⟹ this basket contributes to $\binom{d_i}{t}$ itemsets of size $t$ ⟹ total contribution of $n$ nodes on the right is $\sum_i \binom{d_i}{t}$.

Average of $d_i$ is $d$ so this sum is minimized when each $d_i$ is $d$

⟹ we shall assume that all nodes have the average degree of $d$.

⟹ total contribution of $n$ nodes on the

Right to the counts of itemsets size $t$ is

$$n\left(\frac{d}{t}\right)$$

\# itemsets size $t$ is $\left(\frac{n}{t}\right)$

$\Rightarrow$ Average count of an itemset of size $t$ is $n\left(\frac{d}{t}\right)/\left(\frac{n}{t}\right)$, this expression must

be at least $s$ if we are to argue that $K_{s,t}$ exists

1) FOR $n = 20$, $d = 5$

$\Rightarrow$ $s \le 20\left(\frac{5}{t}\right)/\left(\frac{20}{t}\right)$ $\leftarrow$ guaranted constrain

• if $t = 1$ $\Rightarrow$ $s \le 5$ $\Rightarrow$ $s$ can be as large as $5$.

Note: FOR this example, we can reconfirm $s = 5$ is maximal by consider the following $G$ where node $i^{th}$ on the left is connect to $5$ node $\left(i^{th}; (i+1)^{th}; (i+2)^{th}; (i+3)^{it}; (i+4)^{th}\right)$ on the Right ( if $(i+J)^{th} > 20$ we take $(i+J) \bmod 20$ ). we can see that $G$ does not have $K_{6,1}$.

• if $t \ge 2$ $\Rightarrow$ $s \le 1,053 < t$ ( not considered)

Thus, our $(t,s)$ set is $\boxed{\{(1,5)\}}$

2) FOR $n = 200$, $d = 150$

$$\Rightarrow s \leq 200 \binom{150}{t} / \binom{200}{t} \quad \xleftarrow{\text{guaranteed}} \text{constrain}$$

- $t = 1 \Rightarrow s \leq 150 \Rightarrow S_{max} = 150$ (same arguement in part (a). 1 )

- $t = 2 \Rightarrow s \leq 112.3$ . However if the average support for an item set of 2 is 112.3 , then it is impossible that all those itemsets have support $\leq 112$ . Thus, we can be sure that at least one itemset of size 2 has support 113 or more $\Rightarrow K_{113,2}$ exists (*)

$\Rightarrow S_{max} = 113$ .

- $t = 3 \Rightarrow s \leq 83.95 \Rightarrow S_{max} = 84$ (same arguement with (*)) .

- $t = 4 \Rightarrow s \leq 62.64 \Rightarrow S_{max} = 63$ .

- $t = 5 \Rightarrow s \leq 46.66 \Rightarrow S_{max} = 47$

- $t = 6 \Rightarrow s \leq 34.69 \Rightarrow S_{max} = 35$

- $t = 7 \Rightarrow s \leq 25.75 \Rightarrow S_{max} = 26$

- $t = 8 \Rightarrow s \leq 19.08 \Rightarrow S_{max} = 20$

- $t = 9 \Rightarrow s \leq 14.11 \Rightarrow S_{max} = 15$

- $t = 10 \Rightarrow s \leq 10.41 \Rightarrow S_{max} = 11$

• $t \geq 11 \Rightarrow s \leq 7.67 < t \Rightarrow$ not considered

Thus, our $(s, t)$ maximal set is
$$\{ (1, 150) ; (2, 113) ; (3, 84) ; (4, 63) ; (5, 47)$$
$$(6, 35) ; (7, 26) ; (8, 20) ; (9, 15) ; (10, 11) \}$$

\* Exercise 10.5.2.
Supposed communities $C, D$ with associated probabilities $P_C$ and $P_D$.
a) $C = \{w, x\}$ ; $D = \{y, z\}$.
$P_{wx} = P_C$ ; $P_{yz} = P_D$

$P_{wy} = \epsilon_1$ ; $P_{wz} = \epsilon_2$ ; $P_{xy} = \epsilon_3$ ; $P_{xz} = \epsilon_4$

The likelihood of the graph is :
$$L = P_C \cdot \epsilon_1 \cdot \epsilon_3 \cdot P_D (1 - \epsilon_2)(1 - \epsilon_4)$$

Note that, $\epsilon_i$ is very small so
$$L \approx P_C \, \epsilon_1 \, \epsilon_3 \, P_D \leq \boxed{\epsilon_1 \epsilon_3}. \leftarrow \text{MLE}$$
(a very small number)

$\Rightarrow$ MLE when $P_C = P_D = 1$

b) $C = \{w, x, y, z\}$ ; $D = \{x, y, z\}$

$P_{wx} = P_{wy} = P_{wz} = P_C$

$P_{xy} = P_{yz} = P_{xz} = 1 - (1-P_C)(1-P_D)$

$\qquad\qquad\qquad = P_C + P_D - P_C P_D$

The likelihood of the graph is:

$L = P_C^2 (P_C + P_D - P_C P_D)^2 (1-P_C)^2 (1-P_D)$

let $\quad x = 1-P_C$ ; $y = 1-P_D$ $\quad (x, y \in [0,1])$

$\Rightarrow L = (1-x)^2 (1 - xy)^2 x^2 y$

$= 16 \left[\dfrac{(1-x)}{2} \cdot \dfrac{(1-x)}{2} \cdot x\right] \left[\dfrac{(1-xy)}{2} \cdot \dfrac{(1-xy)}{2} \cdot xy\right]$

$\leq 16 \left(\dfrac{\frac{1-x}{2} + \frac{1-x}{2} + x}{3}\right)^3 \cdot \left(\dfrac{\frac{1-xy}{2} + \frac{1-xy}{2} + xy}{3}\right)^3$

(Cauchy inequality)

$= 16 \left(\dfrac{1}{3}\right)^3 \left(\dfrac{1}{3}\right)^3 = \boxed{\dfrac{16}{729}} \leftarrow$ MLE

Equality holds when : $\begin{cases} \dfrac{1-x}{2} = x \\ \dfrac{1-xy}{2} = xy \end{cases}$

$(\Leftrightarrow) \begin{cases} x = \dfrac{1}{3} \\ y = 1 \end{cases} \quad (\Leftrightarrow) \begin{cases} P_C = \dfrac{2}{3} \\ P_D = 0 \end{cases}$

**Answer to Problem 3**

③. Large-Scale Machine Learning

ⓐ Exercise 12.5.3

a) Gini impurity :

$$G = f(x,y) = 1 - x^2 - y^2$$

(where $x$ is the fraction of examples in 1st class

$y$ " " " 2nd " )

$$\Rightarrow \quad G = 1 - x^2 - (1-x)^2$$

$$= 1 - x^2 - 1 + 2x - x^2$$

$$= 2x - 2x^2$$

$$\Rightarrow G'' = -4 < 0$$

$$\Rightarrow G \text{ is concave}$$

b) Entropy measure :

$$E = g(x,y) = x \log_2\left(\frac{1}{x}\right) + y \log_2\left(\frac{1}{y}\right)$$

$$= x \log_2\left(\frac{1}{x}\right) + (1-x) \log_2\left(\frac{1}{1-x}\right)$$

$$E' = \log_2\left(\frac{1}{x}\right) + x \cdot \frac{-1}{x^2} \cdot (\ln 2 \, x)$$

$$- \log_2\left(\frac{1}{1-x}\right) + (1-x) \cdot \frac{1}{(1-x)^2} \cdot \ln 2 \,(1-x)$$

$$= \log_2\left(\frac{1}{x}\right) - \log_2\left(\frac{1}{1-x}\right)$$

$$\Rightarrow E'' = -\frac{1}{x^2}(\ln 2\, x) - \frac{1}{(1-x)^2}\ln 2\,(1-x)$$

$$= -\ln 2\left(\frac{1}{x} + \frac{1}{1-x}\right) < 0$$

$$\text{(since } 0 < x < 1)$$

$$\Rightarrow \quad E \text{ is concave}$$