S Deepak Narayanan
16110142

4)

a) No. This is because it could so happen that we have a sample like [1] [1] [1] in the first case, whereas that wont occur in the second case (second case).
In second case, we will not have reservoirs holding same values multiple times.
In case 1, it is always the case that could be sampled multiple times.
We have same items sampled multiple times.
We could have all possible samples that we have when we have a reservoir of size of k and in the multiple reservoir case let the vise versa is not true.

Eg: [1] [1] [1] ⟹ k 1 size storage
k size storage ⟹ k 1 size storage

[1|2|3]
↳ New in case 2

b) By induction hypothesis, assume that and we have i+1 elements in stream.
we have reservoir of size k.
Suppose first i elements follow the fact that $\Pr[\text{element in reservoir}] = \frac{k}{i}$
For the $(i+1)$th item,

Probability of it getting added is $\frac{k}{i+1}\left(\frac{1}{k} \neq - \frac{1}{k}\right)$

$= \frac{k}{i+1}$ ( Since we choose up $\frac{k}{i+1}$ $\boxed{\frac{k}{k}=1}$

$(\frac{1}{k} \to$ Pr of replacing an element ) and dont choose up $1 - \frac{k}{i+1}) \to$ for element

Probability of an existing element staying past time $i$ a ~~new element getting chosen~~ and that it it staying in reservoir at time $i$ is not replaced at time $i+1$.

$\Rightarrow \frac{k}{i} \left( 1 - \left(\frac{k}{i+1}\right)\cdot\left(\frac{1}{k}\right)\right) = \frac{k}{i+1}$

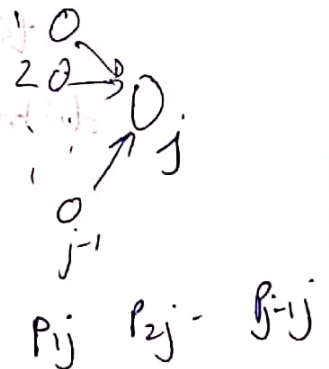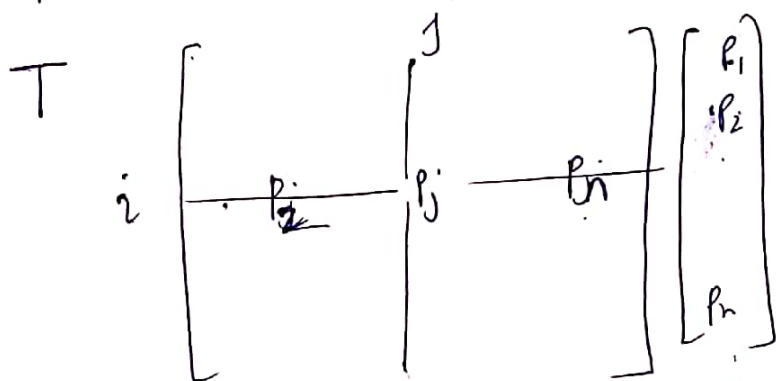(b) Probability an old element stayed in stream at time $i$, and the probability it wasn't kicked out at time $i+1$.

$\Rightarrow \frac{k}{i} \left( 1 - \left(\frac{k}{i+1}\right)\left(\frac{1}{k}\right)\right) = \frac{k}{i+1}$

(c) No. This is not independent sampling. Our sample space reduces if we went to have the following.

Suppose element $x_1$ is present in reservoir and we want to compute probability that $x_2$ is also added to reservoir.

without kicking out $x_1$. ~~This is~~ Assume at time $(t+1)$ $x_2$ comes. If we add $x_2$ to reservoir we need to ensure that $x_1$ is not kicked out,

i.e, Pr of $\left[ x_2 \text{ along with } x_1 \right] = \left( \dfrac{k}{t+1} \right)\left( \dfrac{k-1}{k} \right)$

$\Rightarrow$ Not independent, since this varies as number of samples varies

3) $P \rightarrow$ Probability Vector

$$
T \quad \begin{array}{c} i \\ \\ \end{array}
\begin{bmatrix}
 & & J & & \\
\hline
\cdots & P_i & P_j & P_n & \\
\hline
 & & & & \\
\end{bmatrix}
\begin{bmatrix}
P_1 \\
P_2 \\
\vdots \\
P_n
\end{bmatrix}
$$



$P_{ij} \quad P_{2j} - P_{i-1j}$

$$
\begin{bmatrix}
P_2 & P_3 \\
P_1 & P_2 \\
P_1 & 
\end{bmatrix}
\begin{bmatrix}
P_1 \\
P_2 \\
P_3
\end{bmatrix}
$$

$$
\begin{pmatrix}
P_1 P_2 + P_3 P_3 \\
P_1 P_1 + P_3 P_3 \\
\vdots \\
0
\end{pmatrix}
\begin{pmatrix}
P_2^2 + P_3^2 \\
P_1^2 + P_3^2 \\
0
\end{pmatrix}
\quad P = \quad Tp
$$

$$
= \quad \begin{bmatrix} \\ \\ \\ \end{bmatrix}
$$

for stationarity

$\delta_{ij} = \begin{cases} 1, & \text{edge from } i \rightarrow j \\ 0, & \text{else} \end{cases}$

Added to make stationary

$\delta_{ij} = \begin{cases} \dfrac{1}{\text{edge from } i \rightarrow j} \\ \cancel{\text{else}} \ \text{else} \end{cases}$

$$
P_i = \sum_{j=1}^{n} P_j^2 \, \delta_{ij} + x
$$

$x \Rightarrow$ we can edit only the $p_i^{th}$ element if we have to add

Prc

self loops.     So,                                    $x = P_i x$

$$P_i = \sum_{j=1}^{n} P_j^2 \delta_{ij} + P_i x'$$

$$x' = \frac{1 - \sum_{j=1}^{n} P_j^2 \delta_{ij}}{P_i}$$

We need to add self loops summing upto
. the above value to make $p$, a stationary
distribution :

self loops. So,

$$P_i = \sum_{j=1}^{n} P_j'^2 S_{ij}' + P_i x'$$

$$x = P_i x'$$

$$x' = 1 - \frac{\sum_{j=1}^{n} P_j'^2 S_{ij}'}{P_i}$$

We need to add self loops summing upto the above value to make $P_i$ a stationary distribution.

In this sampling from a random surfer shouldn't affect the efficiency except for the simple fact that we can store the current node and only its neighbours in memory. It is very efficient in this case because we already start off with a stationary distribution and dont need to run random walks a lot of times and because limited many ( walk needs only neighbor data) its efficient.

we have only neighbor data.