

## Quiz 1

Name:

1. Paul has a trouble. He has a big farm and 10,000 chickens. Every day the chickens range free on the farm and come to cages after sunset except exactly one chickens. Every chicken has a number on its wing. The problem is: he wants to know which chicken is missing without using any computer. He is very good at calculation but cannot memorize 9,999 numbers.

How can he find the missing number?

Let  $s = 10,000 * (10,000 + 1) / 2$

With every chicken of ID  $i$ , calculate  $s = s - i$

Finally  $s$  is the missing number!

2. In the farm, Paul has 2 cages. Every day after sunset his chickens go into one of the cages. He wants to know the ratio of common chickens in Cage 1 between yesterday and today. For simplicity, suppose we have 5 chickens {1, 2, 3, 4, 5}. The chickens in Cage 1 yesterday and today are  $S_y = \{1, 3, 5\}$  and  $S_t = \{1, 2, 4, 5\}$  respectively. Approximately calculate the Jaccard coefficient between  $S_y$  and  $S_t$ . You may need some permutations below:

$$P_1 = \langle 4 \ 1 \ 2 \ 5 \ 3 \rangle$$

$$P_2 = \langle 5 \ 3 \ 4 \ 1 \ 2 \rangle$$

$$P_3 = \langle 1 \ 2 \ 5 \ 3 \ 4 \rangle$$

$$P_4 = \langle 3 \ 1 \ 2 \ 5 \ 4 \rangle$$

$$P_5 = \langle 2 \ 5 \ 4 \ 3 \ 1 \rangle$$

Compute signature of  $S_y$  and  $S_t$  respectively.

$$\begin{aligned} \text{Sig}(S_y) = \langle & \min(P_1[1], P_1[3], P_1[5]) = 2, \\ & \min(P_2[1], P_2[3], P_2[5]) = 2, \\ & 1, \\ & 2, \\ & 1 \rangle \end{aligned}$$

Similarly,

$$\text{Sig}(S_t) = \langle 1, 1, 1, 1, 1 \rangle$$

$$\text{Jaccard}(S_y, S_t) = \frac{\# \text{ of common min-hash values}}{\# \text{ of min-hash values}} = \frac{2}{5}$$

Describe a data stream algorithm of computing means and standard deviation with a real number data stream.

```
c = 0
```

```
s = 0
```

```
ss = 0
```

For each element  $v$ ,

```
c += 1
```

```
s += v
```

```
ss += v*v
```

```
mean = s / c
```

```
stddev = sqrt( ss / c - mean * mean)
```

What is difference between Times Series, Cash Register and Turnstile Models?

- \* Time series: simply, keep the incoming element for a state (e.g., sampling)

- \* Cash register: update some states with every incoming element and the states always increase (e.g., Count-min sketching)

- \* Turnstile: update some states with every element and the states sometimes increase, sometimes decrease (e.g., Count sketching)

Suggest an approximate algorithm for the following problem:

Paul sees data stream of positive integers representing  $A_p$  and Carole sees data stream representing  $A_c$ , both on domain  $1, \dots, N$ . Design a streaming algorithm to determine certain number of  $i$ 's with the largest  $\frac{A_p[i]}{A_c[i]}$ .

$S_p = \{\}, S_c = \{\}$

maxratio = -1

for each input  $i$

    Paul computes 10 largest  $A_p[i]$ s using a heap and keeps the corresponding  $i$  in  $S_p$

    Carole computes 10 largest  $A_c[i]$ s using a heap and keep the corresponding  $i$  in  $S_c$

Among common  $i$   $S_p$  and  $S_c$ , answer the largest  $\frac{A_p[i]}{A_c[i]}$

Paul wants to group a data stream of multidimensional data points into  $k$  clusters. Eventually, what he wants compute is  $k$  centers such that the sum of minimum distance from every point to one of them becomes the smallest. That is, the optimal solution is as follows:

Given  $n$  data points  $\{p_1, p_2, \dots, p_n\}$ , the optimal  $k$  centers  $C = \{c_1, c_2, \dots, c_k\}$  of clusters are

$$\operatorname{argmin}_C \sum_{i=1}^n \min_{c \in C} \|p_i - c\|_2^2$$

Provide an approximate algorithm to handle streaming data.

Randomly initialize  $c_j$  ( $j=1, \dots, k$ )

$N_j = 1$  ( $j=1, \dots, k$ )

For each data point  $p_i$

    Find the closest center  $c_j$  among  $c_1, \dots, c_k$

$c_j = (c_j * N_j + p_i) / N_j$

Tell me a problem that requires to handle large streaming data in YOUR research topics (i.e., topics you (or your supervisor) are studying in your lab)