

Task 1

a)

We shall use proof by induction

• Base

Since: $\mu_i^{(0 \rightarrow 1)} = \mu_i^{(0)} + \alpha(x_j' - \mu_i)$ is the update rule, then
 for $n=1$: $\mu_i^{(1)} = \mu_i^{(0)} + \alpha(x_1' - \mu_i^{(0)}) \Rightarrow \mu_i^{(1)} = \mu_i^{(0)} + \alpha x_1' - \alpha \mu_i^{(0)}$
 $\Rightarrow \mu_i^{(1)} = (1-\alpha) \mu_i^{(0)} + \alpha(1-\alpha)^0 x_1'$

• Inductive step

Let's assume that $\mu_i^{(n)} = (1-\alpha)^n \mu_i^{(0)} + \alpha \sum_{k=1}^n (1-\alpha)^{n-k} x_k'$
 holds. We need to prove that the above holds for $n+1$ too.

$$\begin{aligned} & \cancel{\mu_i^{(n)}} \\ & \cancel{(1-\alpha)^n \mu_i^{(0)} + \alpha \sum_{k=1}^n (1-\alpha)^{n-k} x_k'} \end{aligned}$$

$$\begin{aligned} & \text{for } n+1 \text{ the update: } \mu_i^{(n+1)} = \mu_i^{(n)} + \alpha(x_{n+1}' - \mu_i^{(n)}) = \\ & = (1-\alpha)\mu_i^{(n)} + \alpha x_{n+1}' \stackrel{\textcircled{1}}{=} \\ & \Rightarrow \mu_i^{(n+1)} = (1-\alpha)^{n+1} \left[(1-\alpha)^n \mu_i^{(0)} + \alpha \sum_{k=1}^n (1-\alpha)^{n-k} x_k' \right] + \alpha x_{n+1}' \\ & \Rightarrow \mu_i^{(n+1)} = (1-\alpha)^{n+1} \mu_i^{(0)} + \alpha \sum_{k=1}^{n+1} (1-\alpha)^{n+1-k} x_k' = \\ & = \mu_i^{(n+1)} = (1-\alpha)^{n+1} \mu_i^{(0)} + \alpha \sum_{k=1}^{n+1} (1-\alpha)^{n+1-k} x_k' // \quad (\text{proved.}) \end{aligned}$$

b)

$$\mu_i^{(n)} = \frac{1}{n} \sum_{k=1}^n x_k' = \frac{1}{n} \left(\sum_{k=1}^{n-1} x_k' + x_n' \right) = \frac{n-1}{n} \mu_i^{(n-1)} + \frac{1}{n} x_n' =$$

$= \mu_i^{(n-1)} + \frac{1}{n} (x_n' - \mu_i^{(n-1)})$, so the update rule would
 be $\mu_i^{(n)} = \mu_i^{(n-1)} + \frac{1}{n} (x_n' - \mu_i^{(n-1)})$, where n the number of points
 assigned to the cluster up to this point

c) In the first rule, α stays the same. What this does is in essence override (slowly) the contribution of the previous points, which sounds great for online learning.

In the second rule, α the step size slowly decreases, effectively making it so that we don't simply forget the contribution of the previous points but rather consider them all equally. The mean (μ) will converge to the average of the points, making it perfect for any scenario where we know all the points beforehand.

Task 2

a) Let's assume that we have already queried k points (where $k < n$). We now know the true labels of those points, however since we know nothing about the intervals $[a, b]$ and $[c, d]$ or about the distribution of the n points, there is no way we can be sure that our guess for the rest of the $n-k$ points will be correct. (Think of this as a 2 man game where you make the predictions and your opponent chooses a, b, c, d . In this case, you will always lose, unless you query all n points).

b) The idea is the following:

1) Sort the samples

2) Starting from left, query all the samples until you find the first "-1", then skip $m-1$ points by querying like in the beginning

3) Search for the next "-1" and repeat the skipping process

4) Do the same for the next "-1"

5) If you find a "-1" after the last "-1" area then stop, since you know those will all be -1

Overall we perform $n - 3(m-1)$ queries at most

Task 3

a) We know that: $E[R_T] = T\mu^+ - E\left[\sum_{t=1}^T y_t\right] \quad \textcircled{1}$

Rewrite:

$$\begin{aligned} \sum_{t=1}^T y_t &= \sum_{i=1}^k \sum_{j=1}^{n_i^T} y_{i,j} \quad \text{then} \quad E\left[\sum_{i=1}^k \sum_{j=1}^{n_i^T} y_{i,j}\right] = \sum_{i=1}^k E\left[\sum_{j=1}^{n_i^T} y_{i,j}\right] \\ \xrightarrow[\text{tower rule}]{\quad} \sum_{i=1}^k E\left[E\left[\sum_{j=1}^{n_i^T} y_{i,j} | h_i^T\right]\right] &= \sum_{i=1}^k E\left[\sum_{j=1}^{n_i^T} \mu_j\right] = \sum_{i=1}^k E[n_i^T \mu_i] = \end{aligned}$$

$$= \sum_{i=1}^k \mu_i \cdot E[n_i^T]$$

$$\begin{aligned} \textcircled{1} &= T\mu^+ - \sum_{i=1}^k \mu_i E[n_i^T] \quad \cancel{\frac{T\mu^+ - \sum_{i=1}^k \mu_i E[n_i^T]}{T}} = \sum_{i=1}^k E[n_i^T] (\mu^+ - \mu_i) = \\ &= \sum_{i=1}^k E[n_i^T] \Delta_i // \end{aligned}$$

b)

Hoeffding's inequality : $P(|\mu_i - \hat{\mu}_i^+| \geq e) \leq 2 \exp(-2 \frac{e^2}{n_i^+}) \implies$

$$e = \sqrt{\log(T)} / n_i^+ \implies P(|\mu_i - \hat{\mu}_i^+| \geq \sqrt{\frac{\log(T)}{n_i^+}}) \leq 2 \exp(-2 \log(T)) \Rightarrow \\ \Rightarrow P(|\mu_i - \hat{\mu}_i^+| \geq \sqrt{\frac{\log(T)}{n_i^+}}) \leq \frac{2}{T^2} =$$

$$P(|\mu_i - \hat{\mu}_i^+| > \sqrt{\frac{\log(T)}{n_i^+}}) \leq \frac{2}{T^2} \Rightarrow P(\mu_i \notin C_i^+) \leq \frac{2}{T^2}$$

so we ignore " $=$ "

c)

When ~~$\mu_i \in C_i^+$~~ $\mu_i \in C_i^+$: ~~$\hat{\mu}_i^+ \in [\hat{\mu}_i^+ - \sqrt{\frac{\log(T)}{n_i^+}}, \hat{\mu}_i^+ + \sqrt{\frac{\log(T)}{n_i^+}}]$~~

Since $\mu_i \in C_i^+$ and $\mu_i^+ \in C_i^+$: (the above, but also the inverse!)
So: $\hat{\mu}_i^+ \leq \mu_i + \sqrt{\frac{\log(T)}{n_i^+}}$ ② and $\hat{\mu}_i^+ \geq \mu_i - \sqrt{\frac{\log(T)}{n_i^+}}$ (rather than $\mu_i \sim \hat{\mu}_i^+$)

$$\hat{\mu}_i^+ \leq \mu_i^+ + \sqrt{\frac{\log(T)}{n_i^+}} \quad \text{and} \quad \hat{\mu}_i^+ \geq \mu_i^+ - \sqrt{\frac{\log(T)}{n_i^+}} \quad ③$$

also we know that in UCB i would be chosen only if:

$$\hat{\mu}_i^+ + \sqrt{\frac{\log(T)}{n_i^+}} \geq \hat{\mu}_{i^*}^+ + \sqrt{\frac{\log(T)}{n_{i^*}^+}} \quad \textcircled{3}$$

~~from ②, ③, ④~~

we can substitute ② on ③ since and ② on ④ (we will have an even bigger value left with ④ and an even smaller on the right with ②, so it's fine. Then:

$$\hat{\mu}_i^+ + 2\sqrt{\frac{\log(T)}{n_i^+}} \geq \mu^* = 2\sqrt{\frac{\log(T)}{n_{i^*}^+}} > \Delta_i \Rightarrow \frac{4\log(T)}{\Delta_i^2} \geq n_i^+$$

but! accounting for the fact that we try all arms at least once at the start of the algorithm: $n_i^+ \leq \frac{4\log(T)}{\Delta_i^2} + 1$

2)

Task 4

a)

a) Streaming model:

In the streaming model, data can arrive through multiple streams, continuously. The speed at which we intake data is not something we can control, however this distinguishes this model from a database management system, which doesn't need to worry about data passing by too fast and thus being lost.

The incoming data can be stored in memory and thus become accessible later on for queries, however the memory is greatly limited and thus can't hold all the incoming info. A solution to this exists in the form of a huge storage called the archive, however while we can save way more data in the archive, access to it for reads is extremely slow and thus rarely attempted. The last important part of this model are the queries. There are 2 types, standing queries and ad-hoc queries. The former ^{are} comprised of queries that run constantly and can produce answers whenever it's needed. The latter are queries that are asked once about the current state of a stream/streams. The idea

is that you can expect what ~~ad-hoc queries~~ you will receive, either being able to save summaries of the incoming data, or use sliding windows for all streams, which will hold the k most recent data, or data that arrived after a certain point t , so the ad-hoc queries can use this data to answer the needed requests.

b) Counting 1's

The idea is as follows:

* The bucket size starts from 1 always!

You want to count how many 1's approximately showed in the last k bits, $k \leq N$ ($N = \text{window size}$). In that case we will use buckets which will hold 2 pieces of info, the time the last 1 bit was added into the bucket and the size of the bucket. It's important to note that we can only have up to 2 buckets of the same size at any one time, if the number goes to 3 we merge the two oldest buckets and get a bucket with double the size and a time stamp of the ~~newest~~ 1 bit of the two. ~~most recent~~ (in this case we use a time stamp system where the larger the number, the newer it is!) like this:

Lastly we calculate the ~~as~~ ^{like this:} ~~approximate~~

- If $\text{timestamp} - \text{size} \geq \text{current_time} - k$ then count its size as all 1's
- else count half its size

sum for all the buckets and you have a good result!

Example say we are at $t=7$
 $k=5$, 10110110

(also you only end up using $O(\log N)$ space rather than $O(N)$)

$$t=1 \rightarrow [1, 1]$$

$$t=6 \rightarrow [6, 1], [4, 1], [3, 2]$$

$$t=2 \rightarrow [1, 1]$$

$$t=7 \rightarrow [7, 1], [6, 2], [3, 2]$$

$$t=3 \rightarrow [3, 1], [1, 1]$$

full full ↓

$$t=4 \rightarrow [4, 1], [3, 2]$$

half

$$t=5 \rightarrow [4, 1], [3, 2]$$

so we estimate 4 1's, when in reality we have 3, not too bad!

b)

a) BDMO algorithm

This is ~~an algorithm~~ a stream-clustering algorithm that builds upon the counting Δ method we explained and is thus able to maintain clusters efficiently by using summaries (like the number of points in a cluster, the centroid and other info needed to perform merging of clusters). Here the bucket size doesn't need to start from 1.

Algo

Just like before you create the buckets and merge them when 3 buckets of the same size appear.

Inside each bucket you follow a specific clustering strategy (so you can choose to cluster the points that enter) and each time a merge happens you must consider whether or not a merge of clusters should happen, in which case the summary would have to be updated quite a bit.

Finally if a request about the clustering of the last m points is made, then you can be sure that by choosing the covering for the smallest set of buckets ~~that cover~~ ~~the~~ last m points, we can be sure that we will use $\leq 2m$ points, sacrificing some accuracy for the efficiency we presented in the previous question.

b) Parallelism

Given a large set of points ^{for} which we want to compute the centroids, parallelism could be exploited through map-reduce! Each mapper gets a subset of points and ~~yields~~ yields a key-value pair tuple with key = 1 always and value = ~~description~~ of a cluster. Since the keys are all the same, only one reducer

can be used. This reducer simply gets the descriptions of the clusters and tries to merge them, outputting the final clusters.