



## تحلیل داده‌های حجیم

مدرس : دکتر ایمان غلامی

[پاییز ۱۴۰۰]

نگارنده: امیرمحمد شعبانی

تمرین سری ۴: سوال ۱

### الف

با اینکه سوال تعریفی از *don'tknow* نکرده است، اما احتمالاً زمانی این اتفاق می‌افتد که تمام  $k$  سطری که انتخاب کردیم صفر باشند.

طبق فرض سوال، ابتدا یک جایگشت رندم انجام می‌شود و سپس  $k$  سطر را انتخاب می‌کنیم. این برابر است با  $k$  بار انتخاب یک سطر. (چون جایگشت و انتخاب به صورت رندم انجام می‌شود). بنابراین زمانی در این حالت *don'tknow* داریم که تمام سطرهاى انتخابی از سطرهاى صفر باشد. بنابراین چون  $n - m$  سطر غیر صفر داریم، انتخاب یک سطر صفر به احتمال  $\frac{n - m}{n}$  است و چون  $k$  بار انتخاب می‌کنیم؛ پس احتمال آن برابر است با  $(\frac{n - m}{n})^k$ .

### ب

زمانی الگوریتم ناموفق که حداقل دو *don'tknow* داشته باشیم؛ زیرا اگر یکی داشته باشیم، نتیجه می‌شود که با بقیه متفاوت است.

برای محاسبه مقدار  $k$  بنظر چند شرط نیازمندیم که سوال به ما نگفته است. در ابتدا این فرض مهمی است که  $n$  از مقدار  $m$  بسیار بزرگتر باشد. همچنین مقدار  $k$  با مقادیر  $n, m$  بدست می‌آیند.

$$\begin{aligned} \left(\frac{n - m}{n}\right)^k &= \left(1 - \frac{m}{n}\right)^k \Rightarrow \frac{m}{n} \rightarrow 0, \frac{n}{m} = x \Rightarrow \left(1 - \frac{m}{n}\right)^k = \left(1 - \frac{1}{x}\right)^{\frac{k}{x}} = (e^{-1})^{\frac{k}{x}} = e^{-10} \\ \Rightarrow \frac{k}{x} &= 10 \Rightarrow k = 10x \Rightarrow k = \frac{10n}{m} \end{aligned}$$