**I.**

a) Because $X$ follows a discrete probability distribution that describes the probability of $X$ successes (positive cases) in $n$ draws (takes a more extensive PCR test), without replacement, from a finite population (the group that takes basic test) of size $N$ that contains exactly $K$ objects with that feature, where in each draw is either a success (positive) or a failure (negative).

$$P[X = x] = \frac{\binom{K}{x}\binom{N-K}{n-x}}{\binom{N}{n}}, \forall x \in \{\max\{0, n + K - N\}, \dots, \min\{n, K\}\}$$

b) For a population of $N$, we test them one by one, when the $i$-th person is tested positive, i.e. $k_i = 1$, we increment $K$, the count of positive cases within these $N$ people. So $K = \sum_{i=1}^{N} k_i$.

For a population of $N$, we decide one by one if choose them into the smaller group of $n$, when the $i$-th person is tested positive, i.e., $C_i = 1$, we increment $n$, the size of the smaller group within these $N$ people. So $n = \sum_{i=1}^{N} C_i$.

For a population of $N$, we decide one by one if choose them into the smaller group of $n$ and test them, when the $i$-th person is tested positive and he is chosen into smaller group, i.e., $k_i C_i = 1$, we increment $X$, the count of positive cases within the smaller group of size $n$. So $X = \sum_{i=1}^{N} k_i C_i$.

$E[n] = E[\sum_{i=1}^{N} C_i] = \sum_{i=1}^{N} E[C_i] = NE[C_i] = n, \therefore q = E[C_i] = {}^{n}/_{N}$.

$E[K] = E[\sum_{i=1}^{N} k_i] = \sum_{i=1}^{N} E[k_i] = NE[k_i] = K, \therefore E[k_i] = \frac{K}{N}$.

Since $C_i$ and $k_i$ are independent, $E[X] = E[\sum_{i=1}^{N} C_i k_i] = \sum_{i=1}^{N} E[C_i k_i] = NE[C_i]E[k_i] = K, \therefore E[k_i] = np$.

c) $K^2 = (\sum_{i=1}^{N} k_i)(\sum_{j=1}^{N} k_j) = \sum_{i=1}^{N} k_i^2 + \sum_{i\neq j} k_i k_j . \because \forall i \in \{1, \dots, N\}, k_i \in \{0,1\}, \therefore \forall i \in \{1, \dots, N\}, k_i = k_i^2.$

$\therefore K^2 = \sum_{i=1}^{N} k_i + \sum_{i\neq j} k_i k_j = K + \sum_{i\neq j} k_i k_j$

$n^2 = (\sum_{i=1}^{N} C_i)(\sum_{j=1}^{N} C_j) = \sum_{i=1}^{N} C_i^2 + \sum_{i\neq j} C_i C_j . \because \forall i \in \{1, \dots, N\}, C_i \in \{0,1\}, \therefore \forall i \in \{1, \dots, N\}, C_i = C_i^2.$

$\therefore n^2 = \sum_{i=1}^{N} C_i + \sum_{i\neq j} C_i C_j = n + \sum_{i\neq j} C_i C_j$

$X^2 = (\sum_{i=1}^{N} k_i C_i)(\sum_{j=1}^{N} k_j C_j) = \sum_{i=1}^{N} k_i^2 C_i^2 + \sum_{i\neq j} k_i k_j C_i C_j . \because \forall i \in \{1, \dots, N\}, C_i \in \{0,1\},$

$\therefore \forall i \in \{1, \dots, N\}, k_i C_i = k_i^2 C_i^2 . \therefore X^2 = \sum_{i=1}^{N} k_i C_i + \sum_{i\neq j} k_i k_j C_i C_j = X + \sum_{i\neq j} k_i k_j C_i C_j$

$E[X(X - 1)] = \sum_{x=0}^{n} x(x - 1)\frac{\binom{K}{x}\binom{N-K}{n-x}}{\binom{N}{n}} = \sum_{x=2}^{n} \frac{K(K-1)\binom{K-2}{x-2}\binom{N-K}{n-x}}{\frac{N(N-1)}{n(n-1)}\binom{N-2}{n-2}} = \frac{n(n-1)K(K-1)}{N(N-1)}\sum_{x=2}^{n} \frac{\binom{K-2}{x-2}\binom{N-K}{n-x}}{\binom{N-2}{n-2}} = \frac{n(n-1)K(K-1)}{N(N-1)}$

$= \frac{np(n-1)(pN-1)}{(N-1)}$

$E[X^2] = E[X(X - 1)] + E[X] = np + \frac{np(n-1)(pN-1)}{(N-1)}$

$Var(X) = E[X^2] - E[X]^2 = np + \frac{np(n-1)(pN-1)}{N-1} - n^2 p^2 = np\frac{N-1+(n-1)(pN-1)-(N-1)np}{N-1} = \frac{np(1-p)(N-n)}{N-1},$

$Var(X/n) = \frac{Var(X)}{n^2} = \frac{p(1-p)(N-n)}{n(N-1)}$

d) $\because 0 \leq p(1-p) \leq \frac{1}{2} * \frac{1}{2} = \frac{1}{4}, \therefore Var\left(\frac{X}{n}\right) = p(1-p)\frac{N-n}{n(N-1)} \leq \frac{1}{4}\frac{N-n}{n(N-1)}.$

When $N = 1322, n = 1103, Var(X/n) = \frac{K(N-K)(N-n)}{N^2 n(N-1)} = \frac{219K(1322-K)}{2,546,485,692,092}.$

e) According to Chebyshev inequality, $P\left[\left|\frac{X}{n} - p\right| \geq c\right] \leq \frac{Var\left(\frac{X}{n}\right)}{c^2} = 0.05, \therefore c = \sqrt{20Var\left(\frac{X}{n}\right)}$

f) $P[p \in [A, B]] = P\left[\left|\frac{X}{n} - p\right| \leq c\right] \geq 0.95 \rightarrow A = \frac{X}{n} - c, B = \frac{X}{n} + c.$

$. c = \sqrt{20Var\left(\frac{X}{n}\right)} \approx 0.02705. \therefore A \approx 0.55409, B \approx 0.60819.$

g) We use Markov Inequality on the nonnegative random variable $\left(\frac{X}{n} - p\right)^4 . \left\{\left|\frac{X}{n} - p\right| \geq c\right\} \rightarrow \left\{\left(\frac{X}{n} - p\right)^4 \geq c^4\right\},$
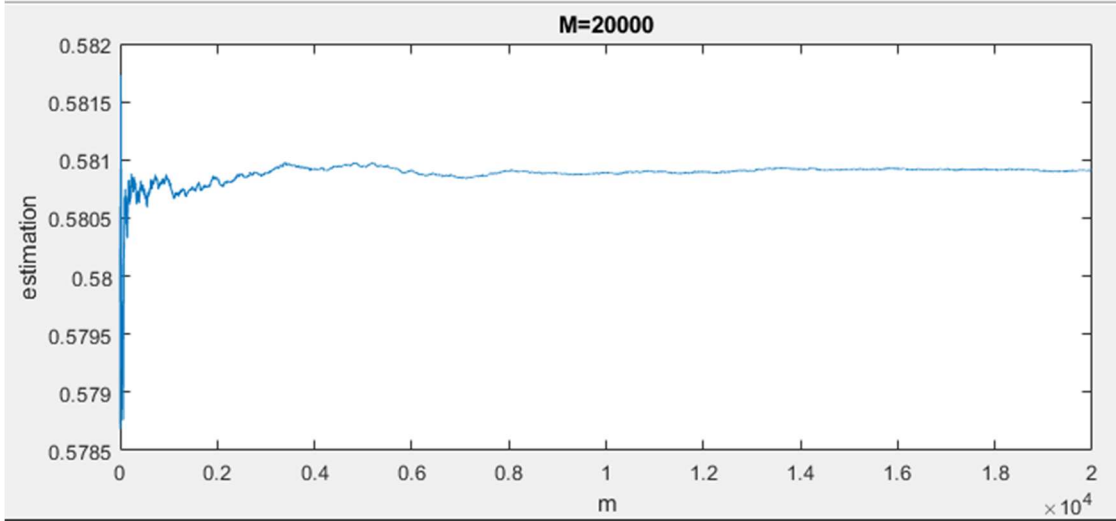
$\therefore P\left[\left|\frac{X}{n} - p\right| \geq c\right] = P\left[\left(\frac{X}{n} - p\right)^4 \geq c^4\right] \leq \frac{E\left[\left(\frac{X}{n}-p\right)^4\right]}{c^4} = \frac{E\left[\left(\frac{1}{n}\right)^4(X-np)^4\right]}{c^4} = \frac{E[(X-np)^4]}{n^4 c^4}.$

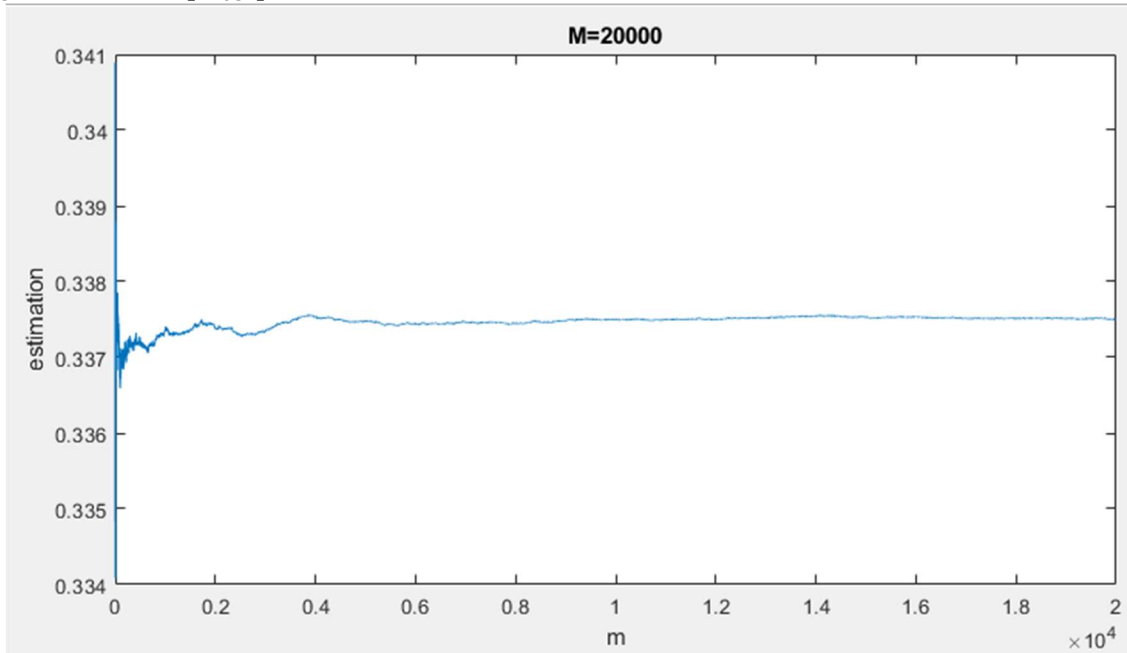$\frac{\mu_4}{\cdot n^4 c^4} = 0.05 \rightarrow c = 0.0168266, A \approx 0.56432, B \approx 0.59797$

h) Chernoff bound: for any $t \in [0, 1-p]$ . $\mathbf{P}[X \geq (p+t)n] \leq e^{-D_{KL}(p+t||p)n}$, where $D_{KL}(p+t||p)$ gives the relative entropy of $p + t$ and $p$.
this paper also provides exponential bounds for hypergeometric distribution with previous studies mentioned.

II

a) Mean of $\frac{1}{j}\sum_{m=1}^{j}\frac{X_m}{1103}$ is 0.5809, variance is 3.1625e-08. $\frac{1}{j}\sum_{m=1}^{j}\frac{X_m}{1103}$ is an unbiased estimator of $p$. We want large $j$ due to Law of Large Numbers.



b) $\frac{1}{j}\sum_{m=1}^{j}\frac{X_m{}^2}{1103^2} = 0.3375; E\left[\frac{X_m{}^2}{1103^2}\right] = 0.3370$



c) For this problem, the upper bound 0.05 is quite loose. Mean of $\frac{1}{j}\sum_{m=1}^{j}1_{\{|\frac{X_m}{1103}-p|\geq c\}}$ is 0.0061, variance is 3.3888e-07.