

## Lecture 3: “Occupancy, Moments and deviations, Randomized selection ”

**Sotiris Nikolettseas**  
**Associate Professor**

CEID - ETY Course  
2013 - 2014

# 1. Some basic inequalities (I)

$$(i) \quad \left(1 + \frac{1}{n}\right)^n \leq e$$

Proof: It is:  $\forall x \geq 0: 1 + x \leq e^x$ . For  $x = \frac{1}{n}$ , we get

$$\left(1 + \frac{1}{n}\right)^n \leq \left(e^{\frac{1}{n}}\right)^n = e$$

$$(ii) \quad \left(1 - \frac{1}{n}\right)^{n-1} \geq \frac{1}{e}$$

Proof: It suffices that  $\left(\frac{n-1}{n}\right)^{n-1} \geq \frac{1}{e} \Leftrightarrow \left(\frac{n}{n-1}\right)^{n-1} \leq e$

But  $\frac{n}{n-1} = 1 + \frac{1}{n-1}$ , so it suffices that  $\left(1 + \frac{1}{n-1}\right)^{n-1} \leq e$   
which is true by (i).

# 1. Some basic inequalities (II)

(iii)  $n! \geq \left(\frac{n}{e}\right)^n$

Proof: It is obviously  $\frac{n^n}{n!} \leq \sum_{i=0}^{\infty} \frac{n^i}{i!}$

But  $\sum_{i=0}^{\infty} \frac{n^i}{i!} = e^n$  from Taylor's expansion of  $f(x) = e^x$ .

(iv) For any  $k \leq n$ :  $\left(\frac{n}{k}\right)^k \leq \binom{n}{k} \leq \left(\frac{ne}{k}\right)^k$

Proof: Indeed,  $k \leq n \Rightarrow \frac{n}{k} \leq \frac{n-1}{k-1}$

Inductively  $k \leq n \Rightarrow \frac{n}{k} \leq \frac{n-i}{k-i}, (1 \leq i \leq k-1)$

Thus  $\left(\frac{n}{k}\right)^k \leq \frac{n}{k} \cdot \frac{n-1}{k-1} \cdots \frac{n-(k-1)}{k-(k-1)} = \frac{n^k}{k!} = \binom{n}{k}$

For the right inequality we obviously have  $\binom{n}{k} \leq \frac{n^k}{k!}$

and by (iii) it is  $k! \geq \left(\frac{k}{e}\right)^k$

## 2. Preliminaries

### (i) Boole's inequality (or union bound)

Let random events  $\mathcal{E}_1, \mathcal{E}_2, \dots, \mathcal{E}_n$ . Then

$$Pr\left\{\bigcup_{i=1}^n \mathcal{E}_i\right\} = Pr\{\mathcal{E}_1 \cup \mathcal{E}_2 \cup \dots \cup \mathcal{E}_n\} \leq \sum_{i=1}^n Pr\{\mathcal{E}_i\}$$

Note: If the events are disjoint, then we get equality.

## 2. Preliminaries

### (ii) Expectation (or Mean)

Let  $X$  a random variable with probability density function (pdf)  $f(x)$ . Its expectation is:

$$\mu_x = E[X] = \sum_x x \cdot Pr\{X = x\}$$

If  $X$  is continuous,  $\mu_x = \int_{-\infty}^{\infty} x f(x) dx$

## 2. Preliminaries

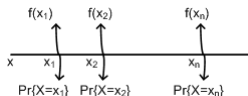
### (ii) Expectation (or Mean)

Properties:

- $\forall X_i \ (i = 1, 2, \dots, n) : E \left[ \sum_{i=1}^n X_i \right] = \sum_{i=1}^n E[X_i]$

This important property is called “linearity of expectation”.

- $E[cX] = cE[X]$ , where  $c$  constant
- if  $X, Y$  stochastically independent, then  
 $E[X \cdot Y] = E[X] \cdot E[Y]$
- Let  $f(X)$  a real-valued function of  $X$ . Then  
$$E[f(x)] = \sum_x f(x) \Pr\{X = x\}$$



## 2. Preliminaries

### (iii) Markov's inequality

Theorem: Let  $X$  a non-negative random variable. Then,  $\forall t > 0$

$$Pr\{X \geq t\} \leq \frac{E[X]}{t}$$

Proof:  $E[X] = \sum_x x Pr\{X = x\} \geq \sum_{x \geq t} x Pr\{X = x\}$

$$\geq \sum_{x \geq t} t Pr\{X = x\} = t \sum_{x \geq t} Pr\{X = x\} = t Pr\{X \geq t\}$$

Note: Markov is a (rather weak) concentration inequality, e.g.

$$Pr\{X \geq 2E[X]\} \leq \frac{1}{2}$$

$$Pr\{X \geq 3E[X]\} \leq \frac{1}{3}$$

etc

## 2. Preliminaries

### (iv) Variance (or second moment)

- Definition:  $Var(X) = E[(X - \mu)^2]$ , where  $\mu = E[X]$   
i.e. it measures (statistically) deviations from mean.
- Properties:
  - $Var(X) = E[X^2] - E^2[X]$
  - $Var(cX) = c^2 Var(X)$ , where  $c$  constant.
  - if  $X, Y$  independent, it is  $Var(X + Y) = Var(X) + Var(Y)$

Note: We call  $\sigma = \sqrt{Var(X)}$  the standard deviation of  $X$ .



## 2. Preliminaries

### (v) Chebyshev's inequality

Theorem: Let  $X$  a r.v. with mean  $\mu = E[X]$ . It is:

$$Pr\{|X - \mu| \geq t\} \leq \frac{Var(X)}{t^2} \quad \forall t > 0$$

Proof:  $Pr\{|X - \mu| \geq t\} = Pr\{(X - \mu)^2 \geq t^2\}$

From Markov's inequality:

$$Pr\{(X - \mu)^2 \geq t^2\} \leq \frac{E[(X - \mu)^2]}{t^2} = \frac{Var(X)}{t^2}$$

Note: Chebyshev's inequality provides stronger (than Markov's) concentration bounds, e.g.

$$Pr\{|X - \mu| \geq 2\sigma\} \leq \frac{1}{4}$$

$$Pr\{|X - \mu| \geq 3\sigma\} \leq \frac{1}{9}$$

etc

### 3. Occupancy - importance

- occupancy procedures are actually stochastic processes (i.e, random processes in time). Particularly, the occupancy process consists in placing randomly balls into bins, one at a time.
- occupancy problems/processes have fundamental importance for the analysis of randomized algorithms, such as for data structures (e.g. hash tables), routing etc.

### 3. Occupancy - definition and basic questions

- general occupancy process: we uniformly randomly and independently put, one at a time,  $m$  distinct objects (“balls”) each one into one of  $n$  distinct classes (“bins”).
- basic questions:
  - what is the maximum number of balls in any bin?
  - how many balls are needed so as no bin remains empty, with high probability?
  - what is the number of empty bins?
  - what is the number of bins with  $k$  balls in them?
- Note: in the next lecture we will study the coupon collector's problem, a variant of occupancy.

### 3. Occupancy - the case $m = n$

Let us randomly place  $m = n$  balls into  $n$  bins.

Question: What is the maximum number of balls in any bin?

Remark: Let us first estimate the expected number of balls in any bin.

For any bin  $i$  ( $1 \leq i \leq n$ ) let  $X_i = \#$  balls in bin  $i$ .

Clearly  $X_i \sim B(m, \frac{1}{n})$  (binomial)

So  $E[X_i] = m \frac{1}{n} = n \frac{1}{n} = 1$

We however expect this “mean” (expected) behaviour to be highly improbable, i.e.,

- some bins get no balls at all
- some bins get many balls

### 3. Occupancy - the case $m = n$

Theorem 1. With probability at least  $1 - \frac{1}{n}$ , no bin gets more than  $k^* = \frac{3 \ln n}{\ln \ln n}$  balls.

Proof: Let  $\mathcal{E}_j(k)$  the event “bin  $j$  gets  $k$  or more balls”. Because of symmetry, we first focus on a given bin (say bin 1). It is  $\Pr\{\text{bin 1 gets exactly } i \text{ balls}\} = \binom{n}{i} \left(\frac{1}{n}\right)^i \left(1 - \frac{1}{n}\right)^{n-i}$  since we have a binomial  $B(n, \frac{1}{n})$ . But

$$\binom{n}{i} \left(\frac{1}{n}\right)^i \left(1 - \frac{1}{n}\right)^{n-i} \leq \binom{n}{i} \left(\frac{1}{n}\right)^i \leq \left(\frac{ne}{i}\right)^i \left(\frac{1}{n}\right)^i = \left(\frac{e}{i}\right)^i$$

(from basic inequality iv)

$$\begin{aligned} \text{Thus } \Pr\{\mathcal{E}_1(k)\} &\leq \sum_{i=k}^n \left(\frac{e}{i}\right)^i \leq \left(\frac{e}{k}\right)^k \cdot \left(1 + \frac{e}{k} + \left(\frac{e}{k}\right)^2 + \dots\right) = \\ &= \left(\frac{e}{k}\right)^k \frac{1}{1 - \frac{e}{k}} \end{aligned}$$

### 3. Occupancy - the case $m = n$

Now, let  $k^* = \lceil \frac{3 \ln n}{\ln \ln n} \rceil$ . Then:

$$Pr\{\mathcal{E}_1(k^*)\} \leq \left(\frac{e}{k^*}\right)^{k^*} \frac{1}{1 - \frac{e}{k^*}} \leq 2 \left(\frac{e}{\frac{3 \ln n}{\ln \ln n}}\right)^{k^*}$$

since it suffices  $\frac{1}{1 - \frac{e}{k^*}} \leq 2 \Leftrightarrow \frac{k^*}{k^* - e} \leq 2 \Leftrightarrow k^* \leq 2k^* - 2e \Leftrightarrow$   
 $\Leftrightarrow k^* \geq 2e$  which is true.

$$\begin{aligned} \text{But } 2 \left(\frac{e}{\frac{3 \ln n}{\ln \ln n}}\right)^{k^*} &= 2 \left(e^{1 - \ln 3 - \ln \ln n + \ln \ln \ln n}\right)^{k^*} \\ &\leq 2 \left(e^{-\ln \ln n + \ln \ln \ln n}\right)^{k^*} \leq 2 \exp\left(-3 \ln n + 6 \ln n \frac{\ln \ln \ln n}{\ln \ln n}\right) \\ &\leq 2 \exp(-3 \ln n + 0.5 \ln n) = 2 \exp(-2.5 \ln n) \leq \frac{1}{n^2} \end{aligned}$$

for  $n$  large enough.

### 3. Occupancy - the case $m = n$

Thus,

$$\begin{aligned} \Pr\{\text{any bin gets more than } k^* \text{ balls}\} &= \Pr\left\{\bigcup_{j=1}^n \mathcal{E}_j(k^*)\right\} \\ &\leq \sum_{j=1}^n \Pr\{\mathcal{E}_j(k^*)\} \leq n\Pr\{\mathcal{E}_1(k^*)\} \leq n\frac{1}{n^2} = \frac{1}{n} \text{ (by symmetry)} \quad \square \end{aligned}$$

### 3. Occupancy - the case $m = n \log n$

- We showed that when  $m = n$  the mean number of balls in any bin is 1, but the maximum can be as high as  $k^* = \frac{3 \ln n}{\ln \ln n}$
- The next theorem shows that when  $m = n \log n$  the maximum number of balls in any bin is more or less the same as the expected number of balls in any bin.
- Theorem 2. When  $m = n \ln n$ , then with probability  $1 - o(1)$  every bin has  $O(\log n)$  balls.



### 3. Occupancy - the case $m = n$ - An improvement

- If at each iteration we randomly pick  $d$  bins and throw the ball into the bin with the smallest number of balls, we can do much better than in Theorem 2:

- Theorem 3. We place  $m = n$  balls sequentially in  $n$  bins as follows:

For each ball,  $d \geq 2$  bins are chosen uniformly at random (and independently). Each ball is placed in the least full of the  $d$  bins (ties broken randomly). When all balls are placed, the maximum load at any bin is at most  $\frac{\ln \ln n}{\ln d} + O(1)$ , with probability at least  $1 - o(1)$  (in other words, a more balanced balls distribution is achieved).

### 3. Occupancy - tightness of Theorem 1

Theorem 1 shows that when  $m = n$  then the maximum load in any bin is  $O\left(\frac{\ln n}{\ln \ln n}\right)$ , with high probability. We now show that this result is tight:

Lemma 1: There is a  $k = \Omega\left(\frac{\ln n}{\ln \ln n}\right)$  such that bin 1 has  $k$  balls with probability at least  $\frac{1}{\sqrt{n}}$ .

Proof:  $Pr[k \text{ balls in bin 1}] = \binom{n}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k}$

$$\geq \binom{n}{k}^k \frac{1}{n^k} \left(1 - \frac{1}{n}\right)^{n-k} \quad (\text{from basic inequality iv})$$
$$= \left(\frac{1}{k}\right)^k \left(1 - \frac{1}{n}\right)^{n-k} \geq \left(\frac{1}{k}\right)^k \left(\frac{1}{2e}\right) = \frac{1}{2e} \left(\frac{1}{k}\right)^k \quad (\text{for } n \geq 2)$$

### 3. Occupancy - tightness of Theorem 1

By putting  $k = \frac{c \ln n}{\ln \ln n}$  we get

$$\Pr\left\{\frac{c \ln n}{\ln \ln n} \text{ balls in bin 1}\right\} \geq \frac{1}{2e} \left(\frac{\ln \ln n}{c \ln n}\right)^{\frac{c \ln n}{\ln \ln n}} \geq \left(\frac{1}{c \ln n}\right)^{\frac{c \ln n}{\ln \ln n}}$$

(for  $n \geq 4$ )

$$= \left(\frac{1}{c 2^{\ln \ln n}}\right)^{\frac{c \ln n}{\ln \ln n}} = \frac{1}{c 2^{\ln \ln n \frac{c \ln n}{\ln \ln n}}} = \frac{1}{c 2^{c \ln n}} = \frac{1}{c n^c} = \Omega(n^{-c})$$

Setting  $c = \frac{1}{2}$  we get  $\Pr\left\{\frac{c \ln n}{\ln \ln n} \text{ balls in bin 1}\right\} \geq \Omega\left(\frac{1}{\sqrt{n}}\right)$   $\square$

### 3. Occupancy - the case $m = n \log n$

Towards a proof of Theorem 2. We use the following bound.

Theorem (Chernoff bound). Let  $X$  a r.v.:

$X = \sum_{i=1}^n X_i = X_1 + \dots + X_n$  where for all  $i$  ( $1 \leq i \leq n$ ) the  $X_i$ 's are independent and

$$X_i = \begin{cases} 1, & \text{with probability } p \\ 0, & \text{with probability } 1 - p \end{cases}$$

Let  $E[X] = np = \mu$ . Then,  $\forall \delta > 0$

$$Pr\{X \geq \mu(1 + \delta)\} \leq \left( \frac{e^\delta}{(1+\delta)^{(1+\delta)}} \right)^\mu \quad \square$$

When placing  $m = n \log n$  balls into  $n$  bins let

$$X_i = \begin{cases} 1, & \text{if ball } i \text{ lands in bin 1 (prob}=\frac{1}{n}) \\ 0, & \text{else} \end{cases}$$

and  $X = \sum_{i=1}^m X_i = \#$  of balls in bin 1. Then  
 $\mu = E[X] = m \frac{1}{n} = \log n$ .

### 3. Occupancy - the case $m = n \log n$

Let us estimate the probability that bin 1 receives more than e.g.  $10 \ln n$  balls

- by the Markov inequality:

$$\Pr\{X \geq 10 \ln n\} \leq \frac{\ln n}{10 \ln n} = \frac{1}{10} \text{ (the bound is not strong)}$$

- by the Chebyshev's inequality:

$X$  is actually binomial, i.e.  $X \sim B(m, \frac{1}{n})$  thus its variance is  $\text{Var}(X) = m \left(\frac{1}{n}\right) \left(1 - \frac{1}{n}\right) = \frac{m}{n} - \frac{m}{n^2} \leq \frac{m}{n}$

$$\text{Thus } \Pr\{X \geq \frac{m}{n} + k\} \leq \Pr\{|X - \frac{m}{n}| \geq k\} \leq \frac{\text{Var}(X)}{k^2} \leq \frac{m}{nk^2}$$

For  $m = n \ln n \Rightarrow \frac{m}{n} = \ln n$  and for  $k = 9 \ln n$  we have

$$\Pr\{X \geq 10 \ln n\} = \Pr\{X \geq \ln n + 9 \ln n\} \leq \frac{n \ln n}{n 81 \ln^2 n} = \frac{1}{81 \ln n}$$

(a bound which is better than the one by Markov's inequality)

### 3. Occupancy - the case $m = n \log n$

Let us estimate the probability that bin 1 receives more than e.g.  $10 \ln n$  balls

- by Chernoff bound:

$$Pr\{X \geq 10 \ln n\} = Pr\{X \geq (1 + 9) \ln n\} \leq \left(\frac{e^9}{10^{10}}\right)^{\ln n} \leq \frac{1}{n^{10}}$$

(much stronger)

Thus,

$$Pr\{\exists \text{ bin with more than } 10 \ln n \text{ balls}\} \leq n \frac{1}{n^{10}} = n^{-9}$$
$$\Rightarrow Pr\{\text{all bins have less than } 10 \ln n \text{ balls}\} \geq 1 - n^{-9}$$

A similar bound applies to the “low tail”, i.e. the probability that there exists a bin with less than, say,  $\frac{1}{10} \ln n$  balls tends to zero, as  $n$  tends to infinity. Overall, there is high concentration around the mean value of  $\ln n$  balls per bin.

### 3. Occupancy - the case $m = n \log n$

Note: The corresponding bounds (for any bin) by Markov's inequality and Chebyshev's inequality are trivial:

- by Markov we get  $\leq \frac{n}{10}$
- by Chebyshev we get  $\leq \frac{n}{81 \ln n}$

### 3. Occupancy - all balls in distinct bins

- Let the experiment of sequentially putting  $m$  balls randomly in  $n$  bins.

Problem: How large  $m$  can be so that the probability of all balls being placed in distinct bins remains high?

- For  $2 \leq i \leq m$ , let  $\mathcal{E}_i =$  “the  $i$ th ball lands in a bin not occupied by the first  $i - 1$  balls”. The desired probability is:

$$\Pr\{\bigcap_{i=2}^m \mathcal{E}_i\} = \prod_{i=2}^m \Pr\{\mathcal{E}_i | \bigcap_{j=2}^{i-1} \mathcal{E}_j\} = \Pr\{\mathcal{E}_2\} \Pr\{\mathcal{E}_3 | \mathcal{E}_2\} \Pr\{\mathcal{E}_4 | \mathcal{E}_2 \mathcal{E}_3\} \cdots \Pr\{\mathcal{E}_m | \mathcal{E}_2 \cdots \mathcal{E}_{m-1}\}$$

$$\text{But } \Pr\{\mathcal{E}_i | \bigcap_{j=2}^{i-1} \mathcal{E}_j\} = 1 - \frac{i-1}{n} \leq e^{-\frac{i-1}{n}}$$

$$\Pr\{\bigcap_{i=2}^m \mathcal{E}_i\} \leq \prod_{i=2}^m e^{-\frac{i-1}{n}} = e^{-\sum_{i=2}^m \frac{i-1}{n}} = e^{-\frac{1}{n} \sum_{i=1}^{m-1} i} = e^{-\frac{m(m-1)}{2n}}$$

Thus, when  $m = \lceil \sqrt{2n} + 1 \rceil$  then this probability is at most  $\frac{1}{e}$  while when  $m$  increases the probability decreases rapidly.

Note: This is similar to the classic “birthday paradox” in probability theory.



## 4. The Randomized Selection Algorithm

- The problem: We are given a set  $S$  of  $n$  distinct elements (e.g. numbers) and we are asked to find the  $k$ th smallest.
- Notation:
  - $r_S(t)$ : the rank of element  $t$  (e.g. the smallest element has rank 1, the largest  $n$  and the  $k$ th smallest has rank  $k$ ).
  - $S_{(i)}$  denotes the  $i$ th smallest element of  $S$  (clearly, we seek  $S_{(k)}$  and  $r_S(S_{(k)}) = k$ ).
- Remark: the fastest known deterministic algorithm needs  $3n$  time and is quite complex. Also, any deterministic algorithm requires  $2n$  time (a tight lower bound).

## 4. The basic idea: random sampling

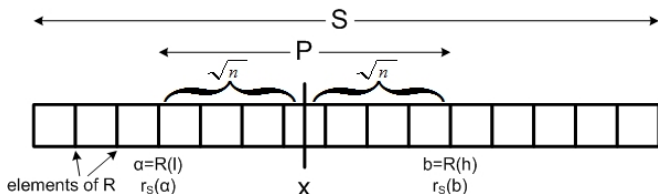
- we will randomly sample a subnet of elements from  $S$ , trying to optimize the following trade-off:
  - the sample should be small enough to be processed (e.g. ordered) in small time
  - the sample should be large enough to contain the  $k$ th smallest element, with high probability

## 4. The Lazy Select Algorithm

- 1 Pick randomly uniformly, with replacement, a subset  $R$  of  $n^{\frac{3}{4}}$  elements from  $S$ .
- 2 Sort  $R$  using an optimal deterministic sorting algorithm.
- 3 Let  $x = k \cdot n^{-\frac{1}{4}}$ .  
 $l = \max\{\lfloor x - \sqrt{n} \rfloor, 1\}$  and  $h = \min\{\lceil x + \sqrt{n} \rceil, n^{\frac{3}{4}}\}$ .  
 $a = R_{(l)}$  and  $b = R_{(h)}$   
By comparing  $a$  and  $b$  to every element of  $S$ , determine  $r_S(a), r_S(b)$ .
- 4 If  $k \in [n^{\frac{1}{4}}, n - n^{\frac{1}{4}}]$ , let  $P = \{y \in S : a \leq y \leq b\}$ .  
Check whether  $S_{(k)} \in P$  and  $|P| \leq 4n^{\frac{3}{4}} + 2$ . If not, repeat steps 1-3 until such a  $P$  is found.
- 5 By sorting  $P$ , identify  $P_{(k-r_S(a)+1)} = S_{(k)}$ .

## 4. Remarks on the Lazy Select Algorithm

- In Step 1, sampling is done with replacement to simplify the analysis. Sampling without replacement is marginally faster but more complex to implement.
- Step 2 takes  $O(n^{\frac{3}{4}} \log n)$  time (which is  $o(n)$ ).
- Step 3 clearly takes  $2n$  time ( $2n$  comparisons). Graphically,



An example: assume  $r_S(a) = 3$  and we want  $S_{(7)}$ . In the sorted list of  $P$  elements,  $S_{(7)} = P_{(k-r_S(a)+1)} = P_{(7-3+1)} = P_5$ , i.e. the 5th element indeed.

## 4. Remarks on the Lazy Select Algorithm

- In Step 4, it is easy to check (in constant time) whether  $S_{(k)} \in P$  by comparing  $k$  to (the now known)  $r_S(a), r_S(b)$ .
- In Step 5, sorting  $P$  takes  $O(n^{\frac{3}{4}} \log n) = o(n)$  time.

Note: we skip in Step 4 the (less interesting) cases where  $k < n^{\frac{1}{4}}$  and  $k > n - n^{\frac{1}{4}}$ . Their analysis is similar.

## 4. When Lazy Select fails?

The algorithm may fail in Step 4, either because  $S_{(k)} \notin P$  because  $|P|$  is large. We will show that the probability of failure is very small.

Lemma 1. The probability that  $S_{(k)} \notin P$  is  $O(n^{-\frac{1}{4}})$ .

Proof: This happens if i)  $S_{(k)} < a$  or ii)  $S_{(k)} > b$ .

i)  $S_{(k)} < a \Leftrightarrow$  fewer than  $l$  ( $l = k \cdot n^{-\frac{1}{4}} - \sqrt{n}$ ) of the samples in  $R$  are less than or equal to  $S_{(k)}$ . Let:

$$X_i = \begin{cases} 1, & \text{the } i\text{th random sample is at most } S_{(k)} \\ 0, & \text{otherwise} \end{cases}$$

Clearly,  $E(X_i) = Pr\{X_i\} = \frac{k}{n}$  and  $Var(X_i) = \frac{k}{n}(1 - \frac{k}{n})$

Let  $X = \sum_{i=1}^{|R|} X_i = \#$  samples in  $R$  that are at most  $S_{(k)}$ . Then

## 4. When Lazy Select fails?

$$\mu_X = E[X] = |R| \cdot E[X_i] = n^{\frac{3}{4}} \frac{k}{n} = kn^{-\frac{1}{4}} \text{ and}$$

$$\sigma_X^2 = Var[X] = \sum_{i=1}^{|R|} Var(X_i) = n^{\frac{3}{4}} \frac{k}{n} \left(1 - \frac{k}{n}\right) \leq \frac{n^{\frac{3}{4}}}{4} \text{ (since the samples are independent)}$$

$$\text{Thus, } Pr\{|X - \mu_X| \geq \sqrt{n}\} \leq \frac{\sigma_X^2}{n} \leq \frac{n^{\frac{3}{4}}}{4n} = O(n^{-\frac{1}{4}})$$

$$\Rightarrow Pr\{X - \mu_X < -\sqrt{n}\} \leq O(n^{-\frac{1}{4}})$$

$$\Rightarrow Pr\{X < \mu_X - \sqrt{n}\} = Pr\{X < \underbrace{kn^{-\frac{1}{4}}}_l - \sqrt{n}\} \leq O(n^{-\frac{1}{4}})$$

## 4. When Lazy Select fails?

ii) The case  $S_{(k)} > b$  is essentially symmetric (at least h of the random samples should be smaller than  $S_{(k)}$ ), so

$$\Pr\{S_{(k)} > b\} = O(n^{-\frac{1}{4}})$$

$$\text{Overall } \Pr\{S_{(k)} \notin P\} = \Pr\{S_{(k)} < a \cup S_{(k)} > b\} = \\ O(n^{-\frac{1}{4}}) + O(n^{-\frac{1}{4}}) = O(n^{-\frac{1}{4}}) \quad \square$$



## 4. The Lazy Select Algorithm

Lemma 2 The probability that  $P$  contains more than  $4n^{\frac{3}{4}} + 2$  elements is  $O(n^{-\frac{1}{4}})$

Proof: Very similar to the proof of Lemma 1: Let

$$k_e = \max\{1, k - 2n^{\frac{3}{4}}\} \text{ and}$$

$$k_n = \min\{k + 2n^{\frac{3}{4}}, n\}$$

If  $S_{(k_l)} < a$  or  $S_{(k_h)} > b$  then  $P$  contains more than  $4n^{\frac{3}{4}} + 2$  elements. For simplicity, let  $k_l = k - 2n^{\frac{3}{4}}, k_h = k + 2n^{\frac{3}{4}}$

Then, it suffices to “simulate” the proof of Lemma 1 for  $k = k_l$  and then for  $k = k_h$ .

## 4. The Lazy Select Algorithm

Theorem The Algorithm Lazy Select finds the correct solution with probability  $1 - O(n^{-\frac{1}{4}})$  performing  $2n + o(n)$  comparisons.

Proof: Due to Lemmata 1, 2 the Algorithm finds  $S_{(k)}$  on the first pass through steps 1-5 with probability  $1 - O(n^{-\frac{1}{4}})$  (i.e., it does not fail in Step 4 avoiding a loop to Step 1). Step 1 obviously takes  $o(n)$  time. Step 2 requires  $O(n^{\frac{3}{4}} \log n) = o(n)$  time, and Step 3 clearly needs  $2n$  comparisons (comparing each of the  $n$  elements of  $S$  to  $a$  and  $b$ ). Overall the time needed is thus  $2n + o(n)$ .