Moments and Deviations

Junhui Jia

{jia@csee.wvu.edu}

For randomized algorithms, usually knowing the bound of expected running time is not enough. It is more desirable to show that the expected running time is small with high probability. To prove this statement, we will begin by examining a family of stochastic processes that is fundamental to the analysis of many types of randomized algorithms. They are Occupancy Problems.

1 Occupancy Problems

A Simple Example of Occupancy Problem

- 1. We have m indistinguishable balls.
- 2. We have n distinct bins.
- 3. We throw those m balls independently, uniformly into those n bins.

Questions:

- 1. What is the **expected number** of balls in a bin?
- 2. What is the **expected number** of bins with k balls in each of it?

Definition: 1.1 The probability that a random variable deviates from its expectation is referred to as the **tail** problem of that deviation.

Theorem: 1.1 The probability of the Union of events is no more than the sum of their probabilities.

$$Pr[\bigcup_{i=1}^{n} E_i] \le \sum_{i=1}^{n} Pr[E_i]$$

<u>Proof</u>: For n = 2

$$P(E_1 \cup E_2) = P(E_1) + P(E_2) - P(E_1E_2) \le P(E_1) + P(E_2)$$

That is

$$Pr[\bigcup_{i=1}^{n=2} E_i] \le \sum_{i=1}^{n=2} Pr[E_i]$$

Suppose for n, it satisfies

$$\Pr[\bigcup_{i=1}^{n} E_i] \le \sum_{i=1}^{n} \Pr[E_i]$$

For n = n + 1, Then we have

$$Pr[\bigcup_{i=1}^{n+1} E_i] = Pr[\bigcup_{i=1}^n E_i \cup E_{n+1}] = Pr[\bigcup_{i=1}^n E_i] + Pr[E_{n+1}] - Pr[\bigcup_{i=1}^n E_i E_{n+1}]$$

Since

 $\Pr[\bigcup_{i=1}^{n} E_i] \le \sum_{i=1}^{n} \Pr[E_i]$

 $Pr[\bigcup_{i=1}^{n} E_i E_{n+1}] \ge 0$

and

Therefore

$$Pr[\bigcup_{i=1}^{n+1} E_i] \le \sum_{i=1}^n Pr[E_i] + Pr[E_{n+1}] \le \sum_{i=1}^{n+1} Pr[E_i]$$

Thereom is proved. \Box

Now, let's consider the case m = n. For $1 \le i \le n$, let X be the number of balls in the ith bin. Let us try to make a statement with very high probability, no bins receives more than k balls, for a chosen k. Let $E_j(k)$ denote the event that bin j has k or more balls in it. The probability that bin j receives exactly i balls is

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

Thus,

$$Pr[E_j(k)] \le \sum_{i=k}^n \left(\frac{e}{i}\right)^i \le \left(\frac{e}{k}\right)^k \left(1 + \frac{e}{k} + \left(\frac{e}{k}\right)^2 + \cdots\right)$$

Let $k^* = \left[(e \lg n) / \lg \lg n \right]$. Substitute k^* into the above equation

$$Pr[E_j(k^*)] \le \left(\frac{e}{k^*}\right)^{k^*} \left(1 + \frac{e}{k^*} + \left(\frac{e}{k^*}\right)^2 + \cdots\right)$$

That is can be simplified as

$$Pr[E_j(k^*)] \le \left(\frac{e}{k^*}\right)^{k^*} \frac{1}{1 - e/k^*}$$

Using $k^* = \lceil (e \lg n) / \lg \lg n \rceil$, we get

$$\Pr[E_j(k^*)] \le n^{-2}$$

We invoke the principle: the probability of the union of the events is no more than their sum. We have

$$Pr[\bigcup_{j=1}^{n} E_j(k^*)] \le \sum_{j=1}^{n} Pr[E_j(k^*)] \le \frac{1}{n}$$

Thus we established

Theorem: 1.2 With probability at least 1 - 1/n, no bin has more than $k^* = (e \ln n) / \ln \ln n$ balls in it.

Suppose that m balls are randomly assigned to n bins. We study the probability of the event that they all land in distinct bins. Consider the assignment of the balls to bins as a sequential process: we throw the first ball into a random bin, then next ball, and so on. For $2 \le i \le m$, let E_i denote the event that the ith ball lands in a bin not containing any of the first i-1 balls. From the probability of the intersection for a collection of event

$$Pr[\bigcap_{i=1}^{k} E_i] = Pr[E_1] \times Pr[E_2|E_1] \times Pr[E_3|E_1 \cap E_2] \cdots Pr[E_k| \cap_{i=1}^{k-1} E_i]$$

We have

$$Pr[\bigcap_{i=2}^{m} E_i] = Pr[E_2]Pr[E_3|E_2]Pr[E_4|E_2 \cap E_3] \cdots Pr[E_m| \cap_{i=2}^{m-1} E_i]$$

The probability that it ball lands in an empty bin given that the first i-1 balls fell into distinct bins is

$$Pr[E_i|\cap_{j=2}^{i-1} E_j] = 1 - \frac{i-1}{n}$$

Making use of the fact that

$$1 - x \le e^{-x}$$

We have

$$\Pr[\bigcap_{i=2}^{m} E_i] \le \prod_{i=2}^{m} \left(1 - \frac{i-1}{n}\right) \le \prod_{i=2}^{m} e^{-(i-1)/n} = e^{-m(m-1)/2n}$$

Corollary: 1.1 For $m = \lceil \sqrt{2n} + 1 \rceil$, the probability that m balls land in distinct bins $\leq \frac{1}{e}$

When m increases beyond this value, the probability drops rapidly. A special case is popular in mathematics as the birthday problem. The 365 days of the year (ignoring leap years) correspond to 365 bins, and the birthday of each of m people is chosen independently and uniformly from 365 days. How large must m be before two people in the same room are likely to share their birthdays? From the Corollary 1.1,

$$m = \left\lceil \sqrt{2n} + 1 \right\rceil = \left\lceil \sqrt{2 \times 365} + 1 \right\rceil = 28$$

Therefore, 28 people is obviously the answer.

2 The Markov and Chebyshev Inequalities

Theorem: 2.1 (Markov Inequality): Let Y be a random variable assuming only non-negative values. Then for all $t \in \mathbb{R}^+$,

$$\Pr[Y \ge t] \le \frac{E[Y]}{t}$$

 $\underline{\operatorname{Proof}}: \ Let$

$$f(y) = 1$$
, if $y \ge t$;
 $f(y) = 0$, otherwise.

For all y, it satisfies

$$f(y) \le \frac{y}{t}$$

Then we have

$$\Pr[Y \ge t] = E[f(Y)] \le E\left[\frac{Y}{t}\right] = \frac{E[Y]}{t}$$

Corollary: 2.1 Let Y be a random variable assuming only non-negative values. Then for all $t \in \mathbb{R}^+$,

$$\Pr[Y \ge kE[Y]] \le \frac{1}{k}$$

Theorem: 2.2 (Chebyshev's Inequality): Let X be a random variable with expectation μ_X and standard deviation σ_X . Then for all $t \in \mathbb{R}^+$,

$$Pr[|X - \mu_X| \ge t\sigma_X] \le \frac{1}{t^2}$$

 $\frac{\text{Proof:}}{\text{First note that}}$

$$|X - \mu_X| \ge t\sigma_X$$
$$(X - \mu_X)^2 \ge t^2 \sigma_X^2$$

The random variable

$$Y = (X - \mu_X)^2$$

has expectation

$$E(Y) = \sigma_X^2$$

From Markov Inequality, we have

$$\Pr[Y \ge t^2 \sigma_X^2] \le \frac{E[Y]}{t^2 \sigma_X^2} \le \frac{\sigma_X^2}{t^2 \sigma_X^2} = \frac{1}{t^2}$$

3 Randomized Selection

Consider the problem of selecting the k^{th} smallest element in a set S of n element. We assume that the elements of S are distinct. Let $r_S(t)$ denote the rank of an element t (the k^{th} smallest element has rank k) and let $S_{(i)}$ denote the i^{th} smallest element of S. Thus the problem becomes that we seek to identify $S_{(k)}$. LazySelect algorithm is introduced.

Let's establish some important property of independent random variables in order to perform the analysis of LazySelect algorithm.

Definition: 3.1 Set X and Y be two random variables defined on the sample space. The joint distribution of X and Y is given by

$$Pr[x, y] = Pr[X = x, Y = y]$$

Theorem: 3.1 The random variable X and Y are independent if

$$Pr[X = x, Y = y] = Pr[X = x]Pr[Y = y]$$

Theorem: 3.2 If X and Y are the independent random variable, then

$$E[XY] = E[X]E[Y]$$

Theorem: 3.3 Let $X_1, X_2, \dots X_m$ be the independent random variables, and $X = \sum_{i=1}^m X_i$. Then

$$\sigma_X^2 = \sum_{i=1}^m \Sigma_{X_i}^2$$

Algorithm LazySelect:

Input: A set S of n elements, and an integer k in [1, n]. **Output:** The k^{th} smallest element of S, $S_{(k)}$.

- 1: Pick $n^{3/4}$ elements from S, chosen independently and uniformly at random with replacement; call this multiset of elements R.
- 2: Sort R in $O(n^{3/4} \log n)$ steps using any optimal sorting algorithm.
- 3: Let $x = kn^{-1/4}$. For $l = max\{\lfloor x \sqrt{n} \rfloor, 1\}$ and $h = min\{\lceil x + \sqrt{n} \rceil, n^{3/4}\}$, let $a = R_{(l)}$ and $b = R_{(h)}$. By comparing a and b to every element of S, determine $r_S(a)$ and $r_S(b)$.
- 4: if k < n^{1/4}, then P = {y ∈ S|y ≤ b}; else if k > n - n^{1/4}, let P = {y ∈ S|y ≥ a}; else if k ∈ [n^{1/4}, n - n^{1/4}], let P = {y ∈ S|a ≤ y ≤ b}; Check whether S_(k) ∈ P and |P| ≤ 4n^{3/4} + 2. If not, repeat Steps 1 - 3 until such a set P is found.
 5: By sorting P in O(|P| log |P|) steps, identify P_{(k-r_S(a)+1)}, which is S_(k).

Algorithm 3.1: LazySelect Algorithm

Proof:

Let μ_i denote $E[X_i]$, and $\mu = \sum_{i=1}^m \mu_i$. The variance of X is given by

$$E[(X - \mu)^2] = E[(\sum_{i=1}^m (X_i - \mu_i))^2]$$

Expanding the latter and using linearity of expectations, we obtain

$$E[(X - \mu)^2] = \sum_{i=1}^m E[(X_i - \mu_i)^2] + 2\sum_{i < j} E[(X_i - \mu_i)(X_j - \mu_j)]$$

Since all pairs of X_i , and X_j are independent, so are the pairs $(X_i - \mu_i)$, $(X_j - \mu_j)$. Each term in the latter summation can be replaced by $E[(X_i - \mu_i)]E[(X_j - \mu_j)]$. Since $E[(X_i - \mu_i)] = E[X_i] - \mu_i = 0$, the latter summation vanishes. It follows that

$$E[(X - \mu)^2] = \sum_{i=1}^m E[(X_i - \mu_i)^2] = \sum_{i=1}^m \sigma_{X_i}^2$$

Thus the idea of the algorithm is to identify two elements a and b in S such that both of the following statements hold with high probability:

- 1. The element $S_{(k)}$ that we seek is in P.
- 2. The set P of elements between a and b is not very large, so that we can sort P inexpensively in step 5.

Theorem: 3.4 With probability $1 - O(n^{-1/4})$, LazySelect finds $S_{(k)}$ on the first pass through Steps 1-5. The running time of LazySelect algorithm is 2n + o(n).

Proof:

The time bound is easily established by examining the algorithm; Step 3 requires 2n comparisons, and all other steps perform o(n) comparisons, provided the algorithm finds $S_{(k)}$ on the first pass through Steps 1–5. We measure the running time of LazySelect algorithm in terms of the number of comparisons performed on it, therefore, the running time of LazySelect algorithm is 2n + o(n).

We now consider the mode of failure: $a > S_{(k)}$ because fewer than l of the samples in R are less than or equal to $S_{(k)}$ (so that $S_{(k)} \notin P$). Set

$$X_i = 1, \ if \ R_{(i)} \le S_{(k)},$$

 $X_i = 0, \ otherwise.$

Thus

$$Pr[X_i = 1] = k/n$$

and

Let

$$X = \sum_{i=1}^{n^3/4} X_i$$

 $Pr[X_i = 0] = 1 - k/n$

be the number of samples of R, that are at most S_k . Note that we really do mean the number of samples, and not the number of distinct elements. The random variables X_i are Bernoulli random variables. Then the expectation and the variance of a Bernoulli random variable with success probability p

$$\mu_X = \frac{kn^{3/4}}{n} = kn^{-1/4}$$
$$\sigma_X^2 = n^{3/4} \left(\frac{k}{n}\right) \left(1 - \frac{k}{n}\right) \le \frac{n^{3/4}}{4}$$
$$l = max\{\lfloor x - \sqrt{n} \rfloor, 1\}$$

Since

Then we have

$$X < x - \sqrt{n}$$
$$X - x < -\sqrt{n}$$
$$|X - x| \ge \sqrt{n}$$

X < l

The probability of the above is

$$Pr[\mid X - x \mid \ge \sqrt{n}] = Pr[\mid X - \mu_X \mid \ge \sqrt{n}]$$

Apply the Chebyshev bound to X and $\sigma_X \leq n^{3/8}/2$

$$Pr[|X - \mu_X| \ge \sqrt{n}] \le Pr[|X - \mu_X| \ge 2n^{1/8}\sigma_X] = O(n^{-1/4})$$

An essentially identical argument shows that

$$Pr[b < S_k] = O(n^{-1/4})$$

Since the probability of the union of events is at most the sum of their probabilities, the probability that either of these events occurs (causing $S_{(k)}$ to lie outside P) is $O(n^{-1/4})$

If the element a is greater than S_k (or if b is smaller than S_k), we fail because P does not contain S_k . Now let's show that

$$Pr[b < S_k] = O(n^{-1/4})$$

Consider the mode of failure: $b < S_{(k)}$ because at least h of the random samples in R should be smaller than $S_{(k)}$ (so that $S_{(k)} \notin P$). Set

$$X_i = 1, \ if \ R_{(i)} \le S_{(k)},$$

 $X_i = 0, \ otherwise.$

 $Pr[X_i = 1] = k/n$

 $Pr[X_i = 0] = 1 - k/n$

Thus

and

Let

$$X = \sum_{i=1}^{n^3/4} X_i$$

be the number of samples of R, that are at most S_k . Note that we really do mean the number of samples, and not the number of distinct elements. The random variables X_i are Bernoulli random variables. Then the expectation and the variance of a Bernoulli random variable with success probability p

$$\mu_X = \frac{kn^{3/4}}{n} = kn^{-1/4}$$

$$_X^2 = n^{3/4} \left(\frac{k}{n}\right) \left(1 - \frac{k}{n}\right) \le \frac{n^{3/4}}{4}$$

$$h = \min\{\lceil x + \sqrt{n}\rceil, n^{3/4}\}$$

$$X \ge h$$

 σ

Since

Then we have

$$X \ge x + \sqrt{n}$$
$$X - x \ge \sqrt{n}$$
$$X - x \ge \sqrt{n}$$

The probability of the above is

$$Pr[|X - x| \ge \sqrt{n}] = Pr[|X - \mu_X| \ge \sqrt{n}]$$

Apply the Chebyshev bound to X and $\sigma_X \leq n^{3/8}/2$

$$Pr[|X - \mu_X| \ge \sqrt{n}] \le Pr[|X - \mu_X| \ge 2n^{1/8}\sigma_X] = O(n^{-1/4})$$

So we proved that

$$Pr[b < S_k] = O(n^{-1/4})$$

The second type of failure occurs when P is too big. To study this, we define $k_l = max\{1, k - 2n^{3/4}\}$ and $k_h = min\{k+2n^{3/4}, n\}$. To obtain an upper bound on the probability of this kind of failure, we will be pessimistic and say that failure occurs if either $a < S_{kl}$ or $b > S_{kh}$. The analysis is very similar to that above in studying the first mode of failure, with k_i and k_h playing the role of k. For $k < n^{1/4}$ and $P = \{y \in S | y \le b\}$, let's show

$$Pr[a < S_{kl}] = O(n^{-1/4})$$

 Set

$$X_i = 1, \ if \ R_{(i)} \le S_{(kl)},$$

 $X_i = 0, \ otherwise.$

Thus

$$Pr[X_i = 1] = k/n$$

and

 $\Pr[X_i = 0] = 1 - k/n$

Let

$$X = \sum_{i=1}^{n^3/4} X_i$$

be the number of samples of R, that are at most S_{kl} . The random variables X_i are Bernoulli random variables. Then the expectation and the variance of a Bernoulli random variable with success probability p

$$\mu_X = \frac{kn^{3/4}}{n} = kn^{-1/4}$$
$$\sigma_X^2 = n^{3/4} \left(\frac{k}{n}\right) \left(1 - \frac{k}{n}\right) \le \frac{n^{3/4}}{4}$$

This implies that $\sigma_X \leq n^{3/8}/2$. Applying the Chebyshev bound to X,

$$Pr[|X - \mu_X| \ge \sqrt{n}] \le Pr[|X - \mu_X| \ge 2n^{1/8}\sigma_X] = O(n^{-1/4})$$

Let's show

$$Pr[b > S_{kh}] = O(n^{-1/4})$$

 $\underline{\text{Proof}}$:

Set

$$X_i = 1, if R_{(i)} \le S_{(kh)},$$

 $X_i = 0, otherwise.$

Thus

$$Pr[X_i = 1] = k/n$$

and

 $\Pr[X_i = 0] = 1 - k/n$

Let

$$X = \sum_{i=1}^{n^3/4} X_i$$

be the number of samples of R, that are at most S_{kh} . The random variables X_i are Bernoulli random variables. Then the expectation and the variance of a Bernoulli random variable with success probability p

$$\mu_X = \frac{kn^{3/4}}{n} = kn^{-1/4}$$
$$\sigma_X^2 = n^{3/4} \left(\frac{k}{n}\right) \left(1 - \frac{k}{n}\right) \le \frac{n^{3/4}}{4}$$

This implies that $\sigma_X \leq n^{3/8}/2$. Applying the Chebyshev bound to X, we have

$$Pr[|X - \mu_X| \ge \sqrt{n}] \le Pr[|X - \mu_X| \ge 2n^{1/8}\sigma_X] = O(n^{-1/4})$$

For $k > n - n^{1/4}$ and $P = \{y \in S | y \ge a\}$; let's show

$$Pr[a < S_{kl}] = O(n^{-1/4})$$

 Set

$$X_i = 1, \ if \ R_{(i)} \le S_{(kl)},$$

 $X_i = 0, \ otherwise.$

 $Pr[X_i = 1] = k/n$

Thus

and

 $Pr[X_i = 0] = 1 - k/n$

Let

$$X = \sum_{i=1}^{n^3/4} X_i$$

be the number of samples of R, that are at most S_{kl} . The random variables X_i are Bernoulli random variables. Then the expectation and the variance of a Bernoulli random variable with success probability p

$$\mu_X = \frac{kn^{3/4}}{n} = kn^{-1/4}$$
$$\sigma_X^2 = n^{3/4} \left(\frac{k}{n}\right) \left(1 - \frac{k}{n}\right) \le \frac{n^{3/4}}{4}$$

This implies that $\sigma_X \leq n^{3/8}/2$. Applying the Chebyshev bound to X,

$$Pr[|X - \mu_X| \ge \sqrt{n}] \le Pr[|X - \mu_X| \ge 2n^{1/8}\sigma_X] = O(n^{-1/4})$$

Let's show

$$Pr[b > S_{kh}] = O(n^{-1/4})$$

 $\underline{\underline{Proof}}$:

Set

$$X_i = 1, if R_{(i)} \le S_{(kh)},$$

 $X_i = 0$, otherwise.

Thus

 $\Pr[X_i = 1] = k/n$

and

$$Pr[X_i = 0] = 1 - k/n$$

Let

$$X = \sum_{i=1}^{n^3/4} X_i$$

be the number of samples of R, that are at most S_{kh} . The random variables X_i are Bernoulli random variables. Then the expectation and the variance of a Bernoulli random variable with success probability p

$$\mu_X = \frac{kn^{3/4}}{n} = kn^{-1/4}$$
$$\sigma_X^2 = n^{3/4} \left(\frac{k}{n}\right) \left(1 - \frac{k}{n}\right) \le \frac{n^{3/4}}{4}$$

This implies that $\sigma_X \leq n^{3/8}/2$. Applying the Chebyshev bound to X, we have

$$Pr[|X - \mu_X| \ge \sqrt{n}] \le Pr[|X - \mu_X| \ge 2n^{1/8}\sigma_X] = O(n^{-1/4})$$

For $k \in [n^{1/4}, n - n^{1/4}]$ and $P = \{y \in S | a \le y \le b\}$; let's show

$$Pr[a < S_{kl}] = O(n^{-1/4})$$

 Set

$$X_i = 1, if R_{(i)} \le S_{(kl)},$$

 $X_i = 0, otherwise.$

Thus

$$Pr[X_i = 1] = k/n$$

and

$$\Pr[X_i = 0] = 1 - k/n$$

Let

$$X = \sum_{i=1}^{n^3/4} X_i$$

be the number of samples of R, that are at most S_{kl} . The random variables X_i are Bernoulli random variables. Then the expectation and the variance of a Bernoulli random variable with success probability p

$$\mu_X = \frac{kn^{3/4}}{n} = kn^{-1/4}$$
$$\sigma_X^2 = n^{3/4} \left(\frac{k}{n}\right) \left(1 - \frac{k}{n}\right) \le \frac{n^{3/4}}{4}$$

This implies that $\sigma_X \leq n^{3/8}/2$. Applying the Chebyshev bound to X,

$$Pr[|X - \mu_X| \ge \sqrt{n}] \le Pr[|X - \mu_X| \ge 2n^{1/8}\sigma_X] = O(n^{-1/4})$$

 \square Let's show

$$Pr[b > S_{kh}] = O(n^{-1/4})$$

 $X_i = 1, if R_{(i)} \le S_{(kh)},$

 $X_i = 0$, otherwise.

 $\underline{\text{Proof}}$:

Set

Thus

$$Pr[X_i = 1] = k/n$$

and

$$\Pr[X_i = 0] = 1 - k/n$$

Let

$$X = \sum_{i=1}^{n^3/4} X_i$$

be the number of samples of R, that are at most S_{kh} . The random variables X_i are Bernoulli random variables. Then the expectation and the variance of a Bernoulli random variable with success probability p

$$\mu_X = \frac{kn^{3/4}}{n} = kn^{-1/4}$$
$$\sigma_X^2 = n^{3/4} \left(\frac{k}{n}\right) \left(1 - \frac{k}{n}\right) \le \frac{n^{3/4}}{4}$$

This implies that $\sigma_X \leq n^{3/8}/2$. Applying the Chebyshev bound to X, we have

$$Pr[|X - \mu_X| \ge \sqrt{n}] \le Pr[|X - \mu_X| \ge 2n^{1/8}\sigma_X] = O(n^{-1/4})$$

Adding up the probability of all of these failure modes, we find that the probability that LazySelect algorithm fail to find a suitable set P is $O(n^{-1/4})$

References

[1] Rajeev Motwani and Prabhakar Raghavan. *Randomized Algorithms*. Cambridge University Press, Cambridge, England, June 1995.