*Dice, used with a plural verb, means small cubes marked with one to six dots, used in gambling games. Dice, used with a singular verb, means a gambling game in which these cubes are used. Dice (plural) can also refer to any small cubes, especially cube-shaped pieces of food (Cut the cheese into dice).*

– MSN Encarta

*Iacta alea est.* (The die is cast.)

–Julius Caesar

*I will never believe that God plays dice with the universe.*

– Albert Einstein

*Thus one comes to perceive, in the concept of independence, at least the first germ of the true nature of problems in probability theory.*

– Kolmogorov

# Lecture VIII
# QUICK PROBABILITY

We review the basic concepts of probability theory, using the axiomatic approach first expounded by A. Kolmogorov. His classic [10] is still an excellent introduction. The axiomatic approach is usually contrasted to the empirical or "Bayesian" approach that seeks to predict real world phenomenon with probabilistic models. Other source books for the axiomatic approach include Feller [5] or the approachable treatment of Chung [4]. Students familiar with probability may use the expository part of this Chapter as reference.

Probability in algorithmics arises in two main ways. In one situation, we have a deterministic algorithm whose input space has some probability distribution. We seek to analyze, say, the expected running time of the algorithm. The other situation is when we have an algorithm that makes random choices, and we analyze its behavior on any input. The first situation is considered less important in algorithmics because we typically do not know the probability distribution on an input space (even if such a distribution exists). By the same token, the second situation derives its usefulness from avoiding any probabilistic assumptions about the input space. Algorithms that make random decisions are said to be **randomized** and comes in two varieties. In one form, the algorithm may make a small error but its running time is worst-case bounded; in another, the algorithm has no error but only its expected running time is bounded. These are known as **Monte Carlo** and **Las Vegas** algorithms, respectively.

There is an understandable psychological barrier to the acceptance of unbounded worst-case running time or errors in randomized algorithms. However, it must be realized that the errors in randomized algorithms are controllable by the user – we can make them as small as we like

at the expense of more computing time. Should we accept an algorithm with error probability of $2^{-99}$? In daily life, we accept and act on information with a much greater uncertainty or likelihood of error than this.

More importantly, randomization is, in many situations, the only effective computational tool available to attack intransigent problems. Until recently, the standard example of a problem not known to be in the class $P$ (of deterministic polynomial time solvable problems), but which admits a randomized polynomial-time algorithm is the **Primality Testing** (deciding if a integer is prime). This problem is now known to be polynomial time, thanks to a breakthrough by M. Agrawal, N. Kayal and N. Saxeena (2002). But the best algorithm has complexity $O(n^{7.5})$ which it is not very practical. Thus randomized primality remains useful in practice. Note that the related problem of factorization of integers does not even have a randomized polynomial time algorithm.

There is a large literature on randomized algorithms. Motwani and Raghavan [12] gives a comprehensive overview of the field. Alon, Spencer and Erdos [1] treats the probabilistic method, not only in algorithmic but also in combinatorial analysis. Karp [7] discusses techniques for the analysis of recurrences that arise in randomized algorithms.

## §1. Axiomatic Probability

All probabilistic phenomena occur in some probabilistic space, which we now formalize (axiomatize).

**¶1. Sample space.** Let $\Omega$ be any non-empty set, possibly infinite. We call $\Omega$ the **sample space** and elements in $\Omega$ are called **sample points**.

We use the following running examples of sample spaces:

(E1) $\Omega = \{H, T\}$ (coin toss). This represents a probabilistic space where there are two outcomes. Typically we identify these outcomes with the results of tossing a coin – head ($H$) or tail ($T$).

(E2) $\Omega = \{1, \ldots, 6\}$ (dice roll). This is a slight variation of (E1) representing the outcomes of the roll of dice, with six possible outcomes.

(E3) $\Omega = \mathbb{N}$ (the natural numbers). This is a significant extension of (E1) since there is a countably infinite number of outcomes.

(E4) $\Omega = \mathbb{R}$ (the real numbers). This is a profound extension because we have gone from discrete space to continuous space.

**¶2. Event space.** Sample spaces becomes more interesting when we give it some structure to form event spaces.

Let $\Sigma \subseteq 2^\Omega$ be a subset of the power set of $\Omega$. The pair $(\Omega, \Sigma)$ is called an **event space** provided three axioms hold:

(A0) $\Omega \in \Sigma$.

(A1)  $A \in \Sigma$ implies that its complement is in $\Sigma$, $\Omega - A \in \Sigma$.

(A2)  If $A_1, A_2, \ldots$ is a countable sequence of sets in $\Sigma$ then $\cup_{i \geq 1} A_i$ is in $\Sigma$.

We call $A \in \Sigma$ an **event** and singleton sets in $\Sigma$ are called **elementary events**. The axioms (A0) and (A1) imply that $\emptyset$ and $\Omega$ are events (the "impossible event" and "inevitable event"). Moreover, the complement of an event is an event.  The set of events are closed under countable unions by axiom (A2). But they are also closed under countable intersections, using de Morgan's law:

*Great, a non-event is an event!*

$$\bigcap_i \overline{A_i} = \overline{\bigcup_i A_i}.$$

In some contexts, an event space may be[1] called a **Borel field** or **sigma field**

We now show two basic ways to construct event spaces:
(i) Construction from generator set: let $\Omega$ be any set and $G \subseteq 2^\Omega$.  Then there is a smallest event space $\overline{G}$ containing $G$; we call this the event space **generated by** $G$. For example, let $\Omega = \{1, \ldots, 6\}$ (dice example, E2). If $G = \{\{1, 2, 3\}, \{3, 4, 5, 6\}\}$ then

*(by the Axiom of Choice)*

$$\overline{G} = \{\emptyset, \{3\}, \{1, 2\}, \{1, 2, 3\}, \{4, 5, 6\}, \{3, 4, 5, 6\}, \{1, 2, 4, 5, 6\}, \{1, 2, 3, 4, 5, 6\}\}.$$

(ii) Subspace construction: if $(\Omega, \Sigma)$ is an event space and $A \in \Sigma$, then we obtain a new event space $(A, \Sigma \cap 2^A)$, which is the **subspace** of $(\Omega, \Sigma)$ **induced** by $A$.

Let $A$ and $B$ be events. We use two standard notations for events: we write "$A^c$" or "$\overline{A}$" for the **complementary event** $\Omega \setminus A$. Also, we write "$AB$" for the event $A \cap B$ (why is this an event?). This is called the **joint event** of $A$ and $B$. Two events $A, B$ are **mutually exclusive** if $AB = \emptyset$.

**¶3. Event Spaces in the Running Examples.**  One choice of $\Sigma$ is

$$\Sigma = 2^\Omega. \tag{1}$$

We call this the **discrete event space**, which is the typical choice for countable sample spaces such as the running examples (E1), (E2) and (E3). In example (E2) the event $\{3, 4, 5, 6\} \in \Sigma$ may be read: "the event of rolling at least 3". What is wrong in assuming (1) in all situations? This choice certainly gives an event space. The problem arises (see next) when we try to assign probabilities to events: for infinite $\Omega$, the choice (1) admits many events for which it is unclear how to assign probabilities. This issue is severe when $\Omega$ is uncountable.

We illustrate the standard way to create an event space for $\Omega = \mathbb{R}$, via a generating set $G \subseteq 2^\Omega$. Let $G$ comprise the half-lines

$$H_r := \{x \in \mathbb{R} : x \leq r\}$$

for each $r \in \mathbb{R}$. This generates an event space $\overline{G}$ that is extremely important.  It is called the **Euclidean Borel field** and denoted $B^1$ or $B^1(\mathbb{R})$. An element of $B^1$ is called an **Euclidean Borel set**. These sets are not easy to describe explicitly, but let us see that some natural sets belongs to $B^1$. We first show that singletons $\{r\}$, $r \in \mathbb{R}$, belong to $B^1$:

*Definition of $B^1 = B^1(\mathbb{R})$*

$$\{r\} = H_r \cap \bigcap_{n \geq 1} H^c_{r+(1/n)}.$$

---

[1] Sometimes, "ring" is used instead of "field".

Then (A2) implies that any countable set belongs to $B^1$. A half-open interval $(a, b] = H_b \setminus H_a = (H_a \cup H_b^c)^c$ belongs to $B^1$. Since singletons are in $B^1$, we now conclude any open or closed interval belongs to $B^1$.

**¶4. Probability space.** So far, we have described concepts that probability theory shares in common with measure theory. Probability properly begins with the next definition: a **probability space** is a triple

*Measure theory is the foundation integral calculus*

$$(\Omega, \Sigma, \Pr)$$

where $(\Omega, \Sigma)$ is an event space and $\Pr : \Sigma \to [0, 1]$ (the unit interval) is a function satisfying the following two axioms:

(P0) $\Pr(\Omega) = 1$.

(P1) if $A_1, A_2, \ldots$ form a countable sequence of pairwise disjoint events then $\Pr(\cup_{i \geq 1} A_i) = \sum_{i \geq 1} \Pr(A_i)$.

We simply call $\Sigma$ the probability space when $\Pr$ is understood. The number $\Pr(A)$ is the **probability** of $A$. A **null event** is one with zero probability. Clearly, the empty set $\emptyset$ is a null event. It is important to realize that when $\Omega$ is infinite, we typically have null events different from $\emptyset$. We deduce that $\Pr(\Omega \setminus A) = 1 - \Pr(A)$ and if $A \subseteq B$ are events then $\Pr(A) \leq \Pr(B)$.

The student should learn to set up the probabilistic space underlying any probabilistic analysis. Whenever there is discussion of probability, you should ask: what is $\Omega$, what is $\Sigma$? This is especially important since probabilists almost never do this explicitly! For finite sample spaces in which each sample point is an event, $\Pr$ is completely specified when we assign probabilities to these (elementary) events.

**¶5. Probability in the Running Examples.** Recall that we have specified $\Sigma$ for each of our running examples (E1)–(E4). We now assign probabilities to events.

In example (E1), we choose $\Pr(H) = p$ for some $0 \leq p \leq 1$. Hence $\Pr(T) = 1 - p$. If $p = 1/2$, we say the coin is **fair**. In example (E2), the probability of an elementary event is $1/6$ in the case of a fair dice.

When $\Omega$ is a finite set, and $\Pr(A) = |A|/|\Omega|$ for all $A \in \Sigma$, we see that the probabilistic framework is simply a convenient language for counting: the size of the set event $A \in \Sigma$ is just its probability times $|\Omega|$—. The space $(\Omega, 2^\Omega, \Pr)$ where $\Pr(\omega) = 1/|\Omega|$ for all $\omega \in \Omega$ is called the **uniform probability model** for $\Omega$.

For (E3), we may choose $\Pr(i) = p_i \geq 0$ ($i \in \Omega = \mathbb{N}$) subject to

$$\sum_{i=0}^{\infty} p_i = 1.$$

An explicit example is illustrated by $p_i = 2^{-(i+1)}$, since $\sum_{i=0}^{\infty} p_i = 2 \sum_{i=0}^{\infty} 2^{-i} = 1$.

For (E4), it is more intricate to define a probability space. We begin with the Euclidean Borel field $B^1 = (\Omega, \Sigma)$, but use the above subspace construction to restrict $B^1$ to a finite

interval $[a, b] \subseteq \mathbb{R}$. The resulting sample space is denoted

$$B^1[a, b] = ([a, b], \Sigma \cap 2^{[a,b]})$$

and it is generated by the intervals $[r, b] = H_r \cap [r, b]$, for all $r \leq c \leq b$. The **uniform probability function** for $B^1[a, b]$ is given by

$$\Pr([r, b]) := (b - r)/(b - a) \tag{2}$$

for all generators $[r, b]$ of $B^1[a, b]$. It is not hard to see that $\Pr(A) = 0$ for every countable $A \in \Sigma$. Thus all countable sets are null events.

It is trickier to define a probability function on $B^1 = B^1(\mathbb{R})$. Thus, there is no analogue of the uniform probability function for $B^1[a, b]$. But a most common way to construct a probability space on $B^1$ is via a continuous function $f : \mathbb{R} \to \mathbb{R}$ with the property that $f(x) \geq 0$, the integral $\int_{-\infty}^{a} f(x)dx$ is defined for all $a \in \mathbb{R}$, and $\int_{-\infty}^{+\infty} f(x)dx = 1$. We call such a function a **density function**. Now define $\Pr(H_a) = \int_{-\infty}^{a} f(x)dx$. For instance, consider the following **Gaussian density function**

$$\phi(x) := \frac{1}{\sqrt{\pi}} e^{-x^2}. \tag{3}$$

Let us verify that it is a density function. Clearly, $\phi(x) \geq 0$. The easiest way to evaluate $\int_{-\infty}^{+\infty} \phi(x)dx$ is to evaluate its square,

$$\left(\int_{-\infty}^{+\infty} \phi(x)dx\right)\left(\int_{-\infty}^{+\infty} \phi(y)dy\right) = \int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty} \phi(x)\phi(y)dxdy$$

$$= \frac{1}{\pi}\int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty} e^{-(x^2+y^2)}dxdy.$$

Next make a change of variable to polar coordinates, $x = r\cos\theta, y = r\sin\theta$. The Jacobian of the change of variable is given by

$$J = \det \frac{\partial(x, y)}{\partial(r, \theta)} = \det \begin{bmatrix} \cos\theta & -r\sin\theta \\ \sin\theta & r\cos\theta \end{bmatrix} = r$$
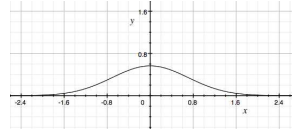
and therefore

$$dxdy = Jdrd\theta = rdrd\theta.$$

We finally get

$$\int_{-\infty}^{+\infty}\int_{-\infty}^{+\infty} e^{-(x^2+y^2)}dxdy = \int_{0}^{2\pi}\int_{0}^{+\infty} e^{-r^2} rdrd\theta$$

$$= \int_{0}^{2\pi} \left[-\frac{e^{-r^2}}{2}\right]_{0}^{\infty} d\theta$$

$$= \int_{0}^{2\pi} \left[\frac{1}{2}\right] d\theta$$

$$= 2\pi \left[\frac{1}{2}\right] = \pi$$

This proves that our original integral is 1. The probability of an arbitrary interval $\Pr[a, b] = \int_{a}^{b} \phi(x)dx$ cannot be expressed as an[2] elementary functions. It can, of course, be numerically approximated to any desired precision. It is also closely related to the error function, $erf(x) := 2\int_{0}^{x} \phi(t)dt$.

---

[2] That is, a univariate complex function obtained from the composition of exponentials, logarithms, nth roots, and the four rational operations $(+, -, \times, \div)$.

**¶6. Constructing Probability Spaces.** We give a product construction for probability spaces: given the event spaces $\Sigma_i \subseteq 2^{\Omega_i}$ $(i = 1, 2)$, let $\Omega := \Omega_1 \times \Omega_2$ and

$$\Sigma_1 \times \Sigma_2 := \{A_1 \times A_2 : A_i \in \Sigma_i, i = 1, 2\}.$$

Note that $A_1 \times A_2$ is empty if $A_1$ or $A_2$ is empty. The event space generated by $\Sigma_1 \times \Sigma_2$ is $(\Omega, \overline{\Sigma_1 \times \Sigma_2})$. If $\mathrm{Pr}_i$ is a probability function for $\Sigma_i$, then we get a probability function $\mathrm{Pr}_1 \times \mathrm{Pr}_2$ for $\overline{\Sigma_1 \times \Sigma_2}$ where

$$(\mathrm{Pr}_1 \times \mathrm{Pr}_2)(A_1 \times A_2) := \mathrm{Pr}_1(A_1)\,\mathrm{Pr}_2(A).$$

We leave it as an exercise to show this $\mathrm{Pr}_1 \times \mathrm{Pr}_2$ can be extended into a probability function for $\overline{\Sigma_1 \times \Sigma_2}$.

We may denote $\Sigma \times \Sigma$ by $\Sigma^2$. For $n \geq 3$, let $\Sigma^n$ denote the event space $(\Sigma^{n-1}) \times \Sigma$. Using this construction, the simple case $\Omega = \{H, T\}$ leads to the event space for a sequence of $n$ coin tosses. We can also let $n = \infty$ and consider $\Sigma^\infty$. This corresponds to sequences of unbounded length: for instance, imagine tossing a coin until we see a head.

**¶7. Decision Tree Models.** An important type of sample space is based on "decision trees". Assuming a finite tree. The sample points $\Omega$ are the leaves of the tree and the sample space is $2^\Omega$. How do we assign probabilities to leaves? At each internal node, we assign a probability of each of its outgoing edges, with the requirement that the sum of these probabilities sum to 1. The uniform probability model satisfies this requirement by assigning each outgoing edge a probability of $1/d$ if the node has degree $d$. Finally, the probability of leaf is just the product of the probability of the edges along the path from the root to the leaf. Intuitively, each sample point amounts to a sequence of probabilistic choices which is made as you move from the root to a leaf.

Let us use a tree model to clarify an interesting puzzle. In a popular TV game show[3] there are three veiled stages. A prize car is placed on a stage behind one of these veils. You are a contestant in this game, and you win the car if you pick the stage containing the car. The rules of the game are as follows: you initially pick one of the stages. Then the game master selects one of the other two stages to be unveiled – this unveiled stage is inevitably car-less. The game master now gives you the chance to switch your original pick. There are two strategies to be analyzed: *always-switch* or *never-switch*. The never-switch strategy is easily analyzed: you have $1/3$ chance of winning. Here are three conflicting claims about the always-switch strategy:

*the game is "Let's Make a Deal"*

CLAIM (I): Your chance of winning is $1/3$, since your chance is the same as the never-switch strategy, as nothing has changed since the start of the game.

CLAIM (II): Your chance of winning is $1/2$, since the car is behind one of the two veiled stages.

CLAIM (III): Your chance of winning is $2/3$, since it is the complement of the never-switch strategy.

Which of these claims (if any) is correct? What is the flaw in at least two of these claims? To solve this puzzle, we set up a tree model for this problem.

The tree model shown in Figure 1 amounts to a sequence of four choices: it alternates between the game master's choice and your choice. First the game master chooses the veiled

---

[3] This problem generated some public interest, including angry letters by professional mathematicians to the New York Times claiming that there ought to be no difference in the two strategies described in the problem.
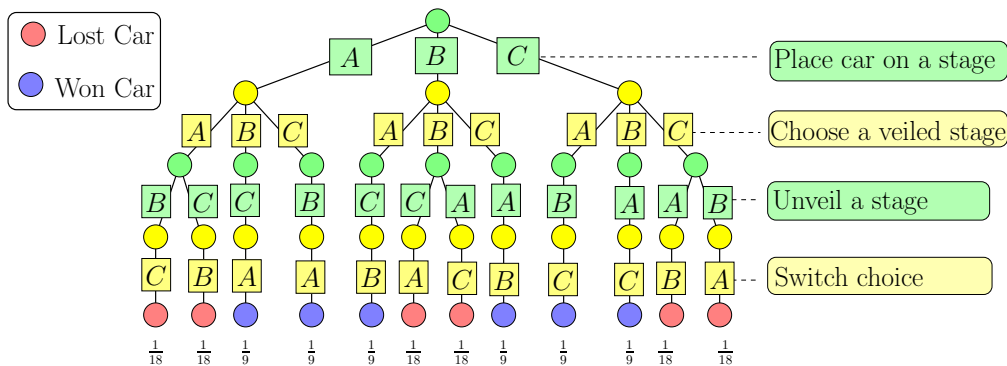
Figure 1: Tree Model for "Let's Make A Deal"

stage to place the car (either $A, B$ or $C$). Then you guess a stage (either $A, B$ or $C$). Then the game master chooses a stage to unveil. Note that if your guess was wrong, the game master has no choice about which stage to unveil. Otherwise, he has two choices. Finally, you choose to switch or not. The figure only shows the always-switch strategy. Using the uniform probability model, the probability of each leaf is either $1/9$ or $1/18$, as shown. Each leaf is colored blue if you win, and colored pink if you lose. Clearly, the probability of winning (blue leaves) is $2/3$. This proves that CLAIM (III) has the correct probability. Indeed, if we use the never-switch strategy, then the pink leaves would be blue and the blue leaves would be pink. This agrees with CLAIM (III) remark that "the always-switch strategy is the complement of the never-switch strategy."

Now we see why CLAIM (I) is wrong is saying "nothing has changed since the start of the game". By revealing one stage, you now know that the car is confined to only two stages. Your original choice has $1/3$ chance of being correct; by switching to the "correct one" among the other two choices, you are actually have $2/3$ chance of being correct. Similarly CLAIM (II) is wrong to say that "the two remaining veiled stages have equal probability of having the car". We agreed that the one you originally chose has $1/3$ chance of having the car. That means the other has $2/3$ chance of having the car.

We can probe a bit further. Do we need the assumption that whenever the game master has a choice of two stages to unveil, both are picked with equal probability? Not at all. In fact the game master can use any rule, including deterministic ones such as choosing to unveil stage $A$ whenever that is possible.

¶8. **Intransitive Spinners.** Consider three spinners, each with three possible outcomes. Spinner $A_0$ has (possible) outcomes $\{1, 6, 8\}$, spinner $A_1$ has outcomes $\{3, 5, 7\}$, and spinner $A_2$ has outcomes $\{2, 4, 9\}$. It is easy to check that the probability that Spinner $A_i$ beats Spinner $A_{i+1}$ is $5/9$ for $i = 0, 1, 2$ (assuming $A_3 = A_0$). So, the relation "spinner $X$ beats spinner $Y$" is not a transitive relation among spinners. This seems to be surprising.

Let us consider the underlying probability spaces. The relation between two spinners is modeled by the space $S \times S = S^2$ where $S = \{1, 2, 3\}$. But suppose we consider all three spinners at once: the space is $S^3 = \{1, 2, 3\}^3$. At any sample point, we do have a total ordering (and hence transitivity) on the three spinners. Thus the binary relation on spinners is a projection from this larger space.

Exercises

**Exercise 1.1:** Show that the method of assigning (uniform) probability to events in $B^1[a, b]$ is well-defined. ◇

**Exercise 1.2:** Let $(\Omega_i, \Sigma_i, \Pr_i)$ be a probability space $(i = 1, 2)$. Recall the product construction for these two spaces.
(a) Show that the event space $\overline{\Sigma_1 \times \Sigma_2}$ is not simply the Cartesian product $\Sigma_1 \times \Sigma_2$. Can you give this a simple description?
(b) Show that the probability function $\Pr = \Pr_1 \times \Pr_2$ is well-defined. ◇

**Exercise 1.3:** Consider the following randomized process, which is a sequence of steps. At each step, we roll a dice that has one of six possible outcomes: $1, 2, 3, 4, 5, 6$. In the $i$-th step, if the outcome is less than $i$, we stop. Otherwise, we go to the next step. For the first step, $i = 1$, and it is impossible to stop at this step. Moreover, we surely stop by the 7-th step. Let $T$ be the random variable corresponding to the number of steps.
(a) Set up the sample space, event space, and probability function for $T$.
(b) Compute the expected value of $T$. ◇

**Exercise 1.4:** The always-switch and never-switch strategies in "Let's Make a Deal" are deterministic. We now consider a randomized strategy: *flip a coin, and to switch only if it is heads*. So this is a mixed strategy – half the time we switch, and half the time we stay put. Analyze this randomized strategy. ◇

**Exercise 1.5:** In "Let's Make a Deal", the game master's probability for initially placing the car, and the probability of unveiling the stage with no car is not under consideration. The game master might even have some deterministic rule for these decisions. Why is this? ◇

**Exercise 1.6:** Let us generalize "Let's Make a Deal". The game begins with a car hidden behind one of $m \geq 4$ possible stages. After you make your choice, the game master unveils all but two stages. Of course, the unveiled stages are all empty, and the two veiled stages always include one you picked.
(a) Analyze the always-switch strategy under the assumption that the game master randomly picks the other stage.
(b) Suppose you want to assume the game master is really trying to work against you. How does your analysis change? ◇

**Exercise 1.7:** J. Quick felt that the sample space $S_n$ we constructed for QuickSort is unnecessarily complicated: why don't we define $S_n$ to be the set of all permutations on the $n$ input numbers. The probability of each permutation in $S_n$ is $1/n!$. What is wrong with this suggestion? ◇

**Exercise 1.8:** (Probabilistic Counters) Recall the counter problem where, given a binary counter $C$ which is initially 0, you can perform the operation `inc(C)` to increments its value by 1. Now we want to do **probabilistic counting**: each time you call `inc(C)`, it will flip a fair coin. If heads, the value of $C$ is incremented and otherwise the value of

$C$ is unchanged. Now, at any moment you could call $\texttt{look}(C)$, which will return *twice* the current value of $C$. Let $X_m$ be the value of $\texttt{look}(C)$ after you have made $m$ calls to $\texttt{inc}(C)$.

(a) Note that $X_m$ is a random variable. What is the sample space $\Omega$ here?

(b) Let $P_m(i)$ be the probability that $\texttt{look}(C) = 2i$ after $m$ $\texttt{inc}$'s. State a recurrence equation for $P_m(i)$ involving $P_{m-1}(i)$ and $P_{m-1}(i-1)$.

(c) Give the exact formula for $P_m(i)$ using binomial coefficients. HINT: you can either use the model in (a) to give a direct answer, or you can try to solve the recurrence of (b). You may recall that binomial identity $\binom{m}{i} = \binom{m-1}{i} + \binom{m-1}{i-1}$.

(d) In probabilistic counting we are interested in the *expected* value of $\texttt{look}(C)$, namely $E[X_m]$. What is the expected value of $X_m$? HINT: express $E[X_m]$ using $P_m(i)$ and do some simple manipulation involving binomial coefficients. If you do not see what is coming out, try small examples like $m = 2, 3$ to see what the answer is.

Remarks: The expected value of $X_m$ can be odd even when the actual value returned is always even. What have we gained by using this counter? We saved 1-bit! But, by a generalization of these ideas, you can probabilistically count to $2^{2^n}$ with an $n$-bit counter, thus saving exponentially many bits. $\diamond$

**Exercise 1.9:** (Intransitive Cubes) Consider 4 cubes whose sides have the following numbers:
$C_0$: $4, 4, 4, 4, 0, 0$
$C_1$: $3, 3, 3, 3, 3, 3$
$C_2$: $6, 6, 2, 2, 2, 2$
$C_3$: $5, 5, 5, 1, 1, 1$
(a) What is the the probability that $C_i$ beats $C_{i+1}$ for all $i = 0, \ldots, 3$? ($C_4 = C_0$)
(b) Let $p \in (0, 1)$ and there exists an assignment of numbers to the faces of the cubes so that the probability that $C_i$ beats $C_{i+1}$ is $p$ for all $i = 0, \ldots, 3$. Show that there exists a biggest value of $p$. E.g., it is easy to see that $p < 1$. Determine this maximum value of $p$. $\diamond$

**Exercise 1.10:** Consider the three spinners $A_0, A_1, A_2$ in the text where $A_i$ beats $A_{i+1}$ with probability $5/9$ ($i = 0, 1, 2$).
(a) Suppose we spin all three spinners at once: what is the probability $p_i$ of $A_i$ beating the others?
(b) Can you design spinners so that the $p_i$'s are equal ($p_0 = p_1 = p_2$) and retain the original property that $A_i$ beats $A_{i+1}$ with a fixed probability $q$? (Originally $q = 5/9$ but you may design a different $q$.) $\diamond$

_____END EXERCISES

## §2. Independence and Conditioning

Intuitively, the outcomes of two tosses of a coin ought to be "independent" of each other. In rolling a pair of dice, the probability of the event "the sum is at least 8" must surely be "conditioned by" the knowledge about the outcome of one of the dice. For instance, knowing that one of the dice is 1 will critically affects this probability. We formalize such ideas of independence and conditioning.

Let $B \in \Sigma$ be any non-null event, *i.e.*, $\Pr(B) > 0$. Such an event $B$ **induces** a probability

*Don't forget that $B$ is non-null in "$\Pr(A|B)$"*

space which we denote by $\Sigma|B$. The sample space of $\Sigma|B$ is $B$ and event space is $\{A \cap B : A \in \Sigma\}$. The probability function $\Pr_B$ of the induced space is given by

$$\Pr_B(A \cap B) = \frac{\Pr(A \cap B)}{\Pr(B)}.$$

It is conventional to write

$$\Pr(A|B)$$

instead of $\Pr_B(A \cap B)$, and to call it the **conditional probability of $A$ given $B$**.

Two events $A, B \in \Sigma$ are **independent** if $\Pr(AB) = \Pr(A)\Pr(B)$.

> For the first time, we have multiplied two probabilities! The product of two probabilities has an interpretation as the probability of the intersection of two corresponding events, but only under an independence assumption. Until now, we have only added probabilities, $\Pr(A) + \Pr(B)$. The sum of two probabilities has an interpretation as the probability of the union of two corresponding events, but only under an disjointness requirement. The combination of adding and multiplying probabilities therefore brings a ring-like structure (involving $+, \times$) into play.

It follows that if $A, B$ are independent then $\Pr(A|B) = \Pr(A)$. More generally, a set $S \subseteq \Sigma$ of events is **independent** if for every $m$-subset $\{B_1, \dots, B_m\} \subseteq S$, we have $\Pr(\cap_{i=1}^m B_i) = \prod_{i=1}^m \Pr(B_i)$. We say $S$ is $k$-**wise independent** if for every $m \le k$, each $m$-subset of $S$ is independent. If $k = 2$, we say $S$ is **pairwise independent**.

Clearly, if $S$ is independent, then it is pairwise independent. The converse may not hold, as shown in this example: Let $\Omega = \{a, b, c, d\}$. With the counting probability model for $\Omega$, consider the events $A = \{a, d\}$, $B = \{b, d\}$, $C = \{c, d\}$. Then we see that $\Pr(A) = \Pr(B) = \Pr(C) = 1/2$ and $\Pr(AB) = \Pr(AC) = \Pr(BC) = 1/4$. Since $\Pr(EF) = \Pr(E)\Pr(F)$ for all events $E \ne F$, we conclude that $S$ is pairwise independent. However $S$ is not independent because $\Pr(ABC) = 1/4 \ne 1/8 = \Pr(A)Pr(B)Pr(C)$.

¶9. **Bayes' Formula.** Suppose $A_1, \dots, A_n$ are mutually exclusive events such that $\Omega = \cup_{i=1}^n A_i$. Then for any event $B$, we have

$$\Pr(B) = \Pr(\uplus_{i=1}^n B \cap A_i) = \sum_{i=1}^n \Pr(B|A_i)\Pr(A_i). \qquad (4)$$

For instance, let $n = 2$ where $A_1$ is "It will rain", $A_2$ is "It is sunny", and $B$ is "I will go to the park". Assume $\Pr(park|rain) = 0$ and $\Pr(park|sunny) = 0.5$. Therefore $\Pr(park) = \Pr(park|rain) + \Pr(park|sunny)$. Thus $\Pr(park) = 0.5$.

Next consider $\Pr(A_j|B) = \Pr(BA_j)/\Pr(B)$. If we replace the numerator by $\Pr(B|A_j)\Pr(A_j)$, and the denominator by (4), we obtain **Bayes' formula** (1763),

$$\Pr(A_j|B) = \frac{\Pr(B|A_j)\Pr(A_j)}{\sum_{i=1}^n \Pr(B|A_i)\Pr(A_i)}. \qquad (5)$$

Continuing our park example, Bayes' formula says

$$\Pr(rain|park) = \frac{\Pr(park|rain)\Pr(rain)}{\Pr(park|rain)\Pr(rain) + \Pr(park|sunny)\Pr(sunny)}.$$

---

Since $\Pr(park|rain) = 0$, we conclude that $\Pr(rain|park) = 0$. In other words, the probability of raining is zero, given that you went to the park. So Bayes' formula "predicted" a sunny day from the fact that I went to the park.

Bayes' formula may be viewed as an inversion of conditional probability: given that $B$ has occurred, you can determine the probability of any of the mutually exclusive $A_j$'s, *assuming you know* $\Pr(B|A_i)$ *for all $i$*. This formula is the starting point for Bayesian probability, the empirical or predictive approach mentioned in the introduction. The goal of Bayesian probability is to use observations to predict the future.

_____EXERCISES

**Exercise 2.1:** (Chung) Recall the Morse code signals of Chapter V. The dot and dash signals are sent in the proportion 3:4. Because of random signal noise, a dot becomes a dash with probability 1/4 (3 becames a 4) and a dash becomes a dot with probability 1/3 (a 4 becomes a 3). Suppose that the probability of sending a dot or a dash is the same.
(a) If a dot is received, what is the probability that it was a dash?
(b) Only four characters in the Morse Code has 2-glyph encodings:

$$\texttt{A: } \cdot - , \quad \texttt{I: } \cdot\cdot, \quad \texttt{M: } --, \quad \texttt{N: } -\cdot$$

So under our dot/dash kind of errors, these four characters can be confused with each other. Compute the probability that the character `A` is sent, given that `A` is received.
◇

_____END EXERCISES

## §3. Random Number Generation and Applications

Randomized algorithms need a source of randomness. Although our main interest in these lectures is randomized algorithms, such a source has many other applications: in the simulation of natural phenomena (computer graphics effects, weather, etc), testing of systems for defects, sampling of populations, decision making and in recreation (dice, card games, etc).

In conventional programming languages such as `Java`, `C`, `C++`, etc, the source of randomness is a function which we will call `random()` (or `rand()`), found in the standard library of that language. Each call to `random()` returns some (presumably) unpredictable a machine-representable floating point number in the half-open interval $[0, 1)$. Mathematically, what we want is a **random number generator** which returns a real number that is uniformly distributed over the unit interval. This distribution is denoted $U_{[0,1]}$. The library function `random()` provides a discrete approximation to $U_{[0,1]}$.

Technically, `random()` is called a **pseudo-random number generator** (PRNG) because it is not truly random. Each PRNG is a parametrized family of number sequences in the unit interval: for each integer parameter $s$, the PRNG defines an infinite sequence

$$X_s = (X_s(1), X_s(2), X_s(3), \ldots)$$

of integers in some range $[0, m)$. This can be interpreted as a sequence over $[0, 1)$ if we divide each $X_s(i)$ by $m$. The $i$-th call to `random()` returns the nearest machine double approximation to $X_s(i)/m$. The sequence $X_s$ is periodic and deterministically generated. The parameter $s$ is freely picked by the user before the initial call to `random()`. We usually call $s$ the **seed** for this particular instantiation of `random()`. In some sense, there is nothing random about $X_s$. On the other hand, the family of $X_s$'s typically satisfy some battery of statistical tests for randomness; this may be sufficient for many applications. The fact that $X_s$ is a deterministic sequence may be important for reproducible in experiments.

¶10. **Linear Congruential Sequences.**  Perhaps the simplest PRNG is the **linear congruential sequences** determined by the choice of three integer parameters: $m$ (the modulus), $a$ (the multiplier) and $c$ (the increment). These parameters, together with seed $s$, determine the sequence

$$X_s(i) = \begin{cases} s & \text{if } i = 0, \\ (aX_s(i-1) + c) \bmod m & \text{if } i \geq 1 \end{cases}.$$

In some PRNG, `random()` returns $E(X_s(i))$ instead of $X_s(i)$ where $E$ is a "bit-extraction routine' that views $X_s(i)$ as a sequence of bits in its binary notation. For example, in `Java`'s `java.util.Random`, we have

*According to Wikipedia...*

$$m = 2^{48}, \quad a = 25,214,903,917, \quad c = 11.$$

Thus, each $X_s(i)$ is an 48-bit integer. The bit-extraction routine of `Java` returns the highest order 32 bits, namely bits 47 to 16:

$$E(X_s(i)) = X_s(i)[47 : 16].$$

On the other hand, the standard library `glibc` used by the compiler `GCC` uses

$$m = 2^{31}, \quad a = 1,103,515,245, \quad c = 12345.$$

Thus each $X_s(i)$ is a 31-bit integer, but the bit extraction routine is a no-op (it returns all 31-bits). The theory of linear congruential sequences provide some easily checkable guidelines for how to choose of $(m, a, c)$.

¶11. **Formula for Joint Events.**  We address another basic "random primitive", random permutations. For instance, this primitive is essential for Implementing a class of random binary trees called treaps (see below). We now show how the preceding `random()` function can be used to compute a random permutation. But first, we need some formulae for computing the probability of a joint event.

From the definition of conditional probability, we have

$$\Pr(AB) = \Pr(A)\Pr(B|A).$$

Then

$$\begin{aligned} \Pr(ABC) &= \Pr(AB)\Pr(C|AB) \\ &= \Pr(A)\Pr(B|A)\Pr(C|AB). \\ \Pr(ABCD) &= \Pr(ABC)\Pr(D|ABC) \\ &= \Pr(A)\Pr(B|A)\Pr(C|AB)\Pr(D|ABC). \end{aligned}$$

More generally,

$$\Pr(A_1 A_2 \cdots A_n) = \prod_{i=1}^{n} \Pr(A_i | A_1 A_2 \cdots A_{i-1}). \tag{6}$$

In proof, simply expand the $i$th factor as $\Pr(A_1 A_2, \ldots, A_i)/\Pr(A_1 A_2, \ldots, A_{i-1})$, and cancel common factors in the numerator and denominator. This formula is "extensible" in that the formula for $\Pr(A_1 \cdots A_n)$ is derived from formula for $\Pr(A_1 \cdots A_{n-1})$ just by appending an extra factor involving $A_n$.

**¶12. Random Permutations.** Fix a natural number $n \geq 2$. Let $S_n$ denote the set of permutations on $[1..n]$. Our problem is to construct a uniformly random element of $S_n$ using a random number generator.

Let us represent a permutation $\pi \in S_n$ by an array $A[1..n]$ where $A[i] = \pi(i)$ $(i = 1, \ldots, n)$. Here is a simple algorithm from Moses and Oakford (see [9, p. 139]).

---

RANDOMPERMUTATION
    Input: an array $A[1..n]$.
    Output: A random permutation of $S_n$ stored in $A[1..n]$.
    1.    for $i = 1$ to $n$ do    ◁ *Initialize array $A$*
    2.        $A[i] = i$.
    3.    for $i = n$ downto 2 do    ◁ *Main Loop*
    4.        $X \leftarrow 1 + \lfloor i \cdot \text{random}() \rfloor$.
    5.        Exchange contents of $A[i]$ and $A[X]$.

---

This algorithm takes linear time; it makes $n - 1$ calls to the random number generator and makes $n - 1$ exchanges of a pair of contents in the array. Here is the correctness assertion for this algorithm:

LEMMA 1. *Every permutation of $[1..n]$ is equally likely to be generated.*

*Proof.* The proof is as simple as the algorithm. Pick any permutation $\sigma$ of $[1..n]$. Let $A'$ be the array $A$ at the end of running this algorithm. It is enough to prove that

$$\Pr\{A' = \sigma\} = \frac{1}{n!}.$$

Let $E_i$ be the event $\{A'[i] = \sigma(i)\}$, for $i = 1, \ldots, n$. Thus

$$\Pr\{A' = \sigma\} = \Pr(E_1 E_2 E_3 \cdots E_{n-1} E_n).$$

First, note that $\Pr(E_n) = 1/n$. Also, $\Pr(E_{n-1}|E_n) = 1/(n-1)$. In general, we see that

$$\Pr(E_i|E_n E_{n-1} \cdots E_{i+1}) = \frac{1}{i}.$$

The lemma now follows from an application of (6) which shows $\Pr(E_1 E_2 E_3 \cdots E_{n-1} E_n) = 1/n!$.
                                                                   **Q.E.D.**

Note that the conclusion of the lemma holds even if we initialize the array $A$ with any permutation of $[1..n]$. This fact is useful if we need to compute another random permutation in the same array $A$.

It is instructive to ask what is the underlying probability space? Basically, if $A'$ is the value of the array at the end of the algorithm, then $A'$ is a random permutation in the sense of §3. That is,

$$A' : \Omega \to S_n$$

where $\Omega$ is a suitable probability space and $S_n$ is the set of $n$-permutations. We can view $\Omega$ as the set $\prod_{i=2}^{n}[0,1)$ where a typical $\omega \in \Omega = (x_2, x_3, \ldots, x_n)$ tells us the sequence of values returned by the $n-1$ calls to the `random()` function.

**Remarks:** Random number generation is an extensively studied topic: Knuth [9] is a basic reference. The concept of randomness is by no means easily pinned down. From the complexity viewpoint, there is a very fruitful approach to randomness called Kolmogorov Complexity. A comprehensive treatment is found in Li and Vitányi [11].

## §4. Random Variables

The concepts so far have not risen much above the level of "gambling and parlor games" (the pedigree of our subject). Probability theory really takes off after we introduce the concept of random variables. Intuitively, a random variable is a function that assigns a real value to each point in a sample space.

As introductory example, consider a discrete event space $(\Omega, 2^\Omega)$ for some finite $\Omega$. Then a random variable is just a function of the form $X : \Omega \to \mathbb{R}$. The event $X^{-1}(c)$ will be written suggestively as "$\{X = c\}$" with probability "$\Pr\{X = c\}$". Clearly, for most $c \in \mathbb{R}$, $\Pr\{X = c\}$ would be 0; Using our running example (E1) where $\Omega = \{H, T\}$, consider the random variable $X$ given by the assignment $X(H) = 100$ and $X(T) = 0$. This random variable might represent the situation where you win \$100 if the coin toss is a head, and you win nothing if a tail. Next suppose the probabilities are given by $\Pr\{X = 1\} = p$ and $\Pr\{X = 0\} = q = 1 - p$. In this case, we may say that your "expected win" is $\Pr\{X = 1\} \cdot X(H) + \Pr\{X = T\} \cdot X(T) = p \cdot 100 + q \cdot 0 = 100p$. I.e., in this game, you expect to win \$100p. As we shall see, a basic thing we do with random variables is to compute their expected values.

In the above example, it was easy to assign a probability to events such as $\{X = c\}$. But if $\Omega$ is infinite, this is trickier. The solution relies on the concept of Borel fields.

We define a **random variable** (r.v.) of a probability space $(\Omega, \Sigma, \Pr)$ is an real function

$$X : \Omega \to \mathbb{R}$$

such that for all $c \in \mathbb{R}$,

$$X^{-1}(H_c) = \{\omega \in \Omega \ : \ X(\omega) \le c\}$$

belongs to $\Sigma$, where $H_c$ is a generator of the Euclidean Borel field $B^1$. Sometimes the range of $X$ is the extended reals $\mathbb{R} \cup \{\pm\infty\}$. It follows that for any Euclidean Borel set $A \in B^1$, the set

$$X^{-1}(A) = \{\omega \in \Omega \ : \ X(\omega) \in A\} \tag{7}$$

is an event. This event is usually written

$$\{X \in A\}. \tag{8}$$

---

**Convention.** Writing (8) for (7) illustrates the habit of probabilists to avoid explicitly mentioning sample points. More generally, probabilists will specify events by writing $\{\ldots X \ldots Y \ldots\}$ where "$\ldots X \ldots Y \ldots$" is some predicate on r.v.'s $X, Y$, etc. This really denotes the event $\{\omega \in \Omega \ : \ \ldots X(\omega) \ldots Y(\omega) \ldots\}$. For instance, $\{X \le 5, X + Y > 3\}$ refers to the event $\{\omega \in \Omega : X(\omega) \le 5, X(\omega) + Y(\omega) > 3\}$. Moreover, instead of writing $\Pr(\{\ldots\})$, we simply write $\Pr\{\ldots\}$, where the pair of curly brackets reminds us that $\{\ldots\}$ is a set (which happens to be an event).

---

An important property about r.v.'s is their closure property under the usual numerical operations: if $X, Y$ are r.v.'s then so are

$$\min(X, Y), \quad \max(X, Y), \quad X + Y, \quad XY, \quad X^Y, \quad X/Y.$$

In the last case, we require $Y(\omega) \neq 0$ for all $\omega \in \Omega$.

**¶13. Two Types of Random Variables.** Practically all random variables in probability theory are either discrete or continuous. Random variables over finite sample spaces are easy to understand, and falls under the discrete case. However, the discrete case covers a bit more.

A r.v. $X$ is **discrete** if the range of $X$,

$$X(\Omega) := \{X(\omega) : \omega \in \Omega\}$$

is countable. Clearly if $\Omega$ is countable, then $X$ is automatically discrete. The simple but important case is when $X(\Omega)$ has size 2; we then call $X$ a **Bernoulli r.v.**. Typically, $X(\Omega) = \{0, 1\}$ as in our introductory example of coin tossing (that is why coin tossing experiments are often called Bernoulli trials). When $X(\Omega) = \{0, 1\}$, we also call $X$ the **indicator function** of the event $\{X = 1\}$. Thus $X$ is the indicator function for the "head event" for the coin tossing example. In discrepancy theory, the Bernoulli r.v. typically have $X(\Omega) = \{+1, -1\}$.

Suppose the implicit probability space is either $B^1$ or $B^1[a, b]$. A r.v. $X$ is **continuous** if there exists a nonnegative function $f(x)$ defined for all $x \in \mathbb{R}$ such that for any Euclidean Borel set $A \in B^1$, the probability is given by an integral:

$$\Pr\{X \in A\} = \int_A f(x)dx$$

(cf. (8)). In particular, for an interval $A = [a, b]$, we have $\Pr\{X \in A\} = \Pr\{a \leq X \leq b\} = \int_a^b f(x)dx$. For instance, this implies $\Pr\{X = a\} = 0$ for all $a$. We call $f(x)$ the **density function** of $X$. It is easy to see that density functions must be non-negative $f(x) \geq 0$ and $\int_{-\infty}^{+\infty} f(x)dx = 1$.

**¶14. Expected Values.** In our introductory example, we noted that a basic application of r.v.'s is to compute their average or expected values. If $X$ is a discrete r.v. whose range is

$$\{a_1, a_2, a_3 \ldots\} \tag{9}$$

then its **expectation** (or, **mean**) $\mathbb{E}[X]$ is defined to be

$$\mathbb{E}[X] := \sum_{i \geq 1} a_i \Pr\{X = a_i\}. \tag{10}$$

This is well-defined provided the series converges absolutely, *i.e.*, $\sum_{i \geq 1} |a_1| \Pr\{X = a_i\}$ converges. If $X$ is a continuous r.v., we must replace the countable sum in (10) by an integral. If the density function of $X$ is $f(x)$ then its expectation is defined as

$$\mathbb{E}[X] := \int_{-\infty}^{\infty} u f(u)du. \tag{11}$$

Note that if $X$ is the indicator variable for an event $A$ then

$$\mathbb{E}[X] = \Pr(A).$$

As examples of random variables, suppose in running example (E1), if we define $X(H) = 1, X(T) = 0$ then $X$ is the indicator function of the "head event". For (E2), let us define $X(i) = i$ for all $i = 1, \ldots, 6$. If we have a game in which a player is paid $i$ dollars whenever the player rolls an outcome of $i$, then $X$ represents "payoff function".

**¶15. How long do I wait for success?** An infinite sequence of Bernoulli r.v.s,

$$X_1, X_2, X_3, \ldots$$

is called a **Bernoulli process**. Each $X_i$ is called a **trial**. *We assume that each trial is independent of the previous ones.* So the sample space is $\Omega = \{0, 1\}^\infty$. Here are two sample points: $\omega_1 = (1, 1, 1, 1, 1, \ldots)$ (all 1's) and $\omega_2 = (0, 0, 1, 0, 0, 1, 0 \ldots)$ (two 0's followed by a 1, repeated in this way). The event space $\Sigma$ is a bit harder to describe (Exercise). Let us suppose $\Pr\{X = 1\} = p$ and $\Pr\{X = 0\} = q := 1 - p$. This defines a probability function $\Pr : \Sigma \to [0, 1]$. For instance, the event $\{X_1 = 1, X_3 = 0\}$ has probability $p(1 - p)$.

Let us interpret the event $\{X = 1\}$ as "success" and the event $\{X = 0\}$ as "failure". Let $Y$ denote the random variable indicating the number of trials until we encounter the first success. For instance, $Y(\omega_1) = 1$ and $Y(\omega_2) = 3$ in the earlier sample points. We can now ask: what is the expected value of $Y$? The answer is very simple: $\mathbf{E}[Y] = 1/p$. For example, if you keep tossing a fair coin, your expect number of tosses before you see a head is $1/p = 2$. Let us prove this.

$$
\begin{aligned}
\mathbf{E}[Y] &= \sum_{i=1}^{\infty} i \cdot \Pr\{Y = i\} \\
&= \sum_{i=1}^{\infty} i q^{i-1} p \\
&= p \sum_{i=1}^{\infty} i q^{i-1} \\
&= p \frac{1}{(1-q)^2} \\
&= p \frac{1}{p^2} = 1/p.
\end{aligned}
$$

## §5. Random Objects

We seek to generalize the concept of random variables. Let $D$ be a set and $(\Omega, \Sigma, \Pr)$ a probability space. A **random function over** $D$ is a function

$$f : \Omega \to D$$

such that for each $x \in D$, the set $f^{-1}(x)$ is an event. Here, $(\Omega, \Sigma, \Pr)$ is the **underlying probability space** of $f$.

We may speak of the **probability** of $x \in D$, *viz.*, $\Pr(f^{-1}(x)) = \Pr\{f = x\}$. We say $f$ is **uniformly distributed on** $D$ if $\Pr(f^{-1}(x)) = \Pr(f^{-1}(y))$ for all $x, y \in D$. We often also use bold fonts ($\mathbf{f}$ instead of $f$, etc) to denote random functions. If the elements of $D$ are objects of some category $T$ of objects, we may call $f$ a **random $T$ object**. If $D$ is a finite set, we also say $f$ is **uniformly random** if $\Pr\{f = d\} = 1/|D|$ for all $d \in D$. Here are some common choices of $D$:

- If $D = \mathbb{R}$, then $f$ is a **random number**. Of course, this is the same concept as "random variable" which we defined earlier.

- If $D$ is some set of graphs we call $f$ a **random graph**. E.g., $D$ is the set of bigraphs on $\{1, 2, \ldots, n\}$.

- For any finite set $S$, we call $f$ a **random $k$-set** of $S$ if $D = \binom{S}{k}$.

- If $D$ is the set of permutations of $S$, then $f$ is a **random permutation** of $S$.

- For any set $D$, we may call $f$ a **random $D$-element**.

The power of random objects is that they are composites of the individual objects of $D$. For many purposes, these objects are as good as the honest-to-goodness objects in $D$. Another view of this phenomenon is to use the philosophical idea of alternative or possible worlds. Each $\omega \in \Omega$ is a possible world. Then $f(\omega)$ is just the particular incarnation of $f$ in the world $\omega$.

*Good, $\omega$ means* *<u>world</u>!*

**¶16. Example:** (Random Points in Finite Fields) Consider the uniform probability space on $\Omega = F^2$ where $F$ is any finite field. The simplest examples of such fields are $\mathbb{Z}_p$ (integers mod $p$ for some prime $p$). For each $x \in F$, consider the random function
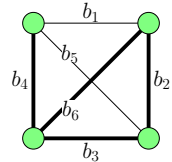
$$\mathbf{h}_x : \Omega \to F,$$
$$\mathbf{h}_x(\langle a, b \rangle) = ax + b, \qquad (\langle a, b \rangle \in \Omega).$$

We claim that $\mathbf{h}_x$ is a uniform random element of $F$, *i.e.*, $\Pr\{\mathbf{h}_x = i\} = 1/|F|$ for each $i \in F$. This amounts to saying that there are exactly $|F|$ sample points $\langle a, b \rangle = \omega$ such that $\mathbf{h}_x(\omega) = i$. To see this, consider two cases: (1) If $x = 0$ then $h_0(\langle a, b \rangle) = i$ implies $b = i$. But $a$ can be any element in $F$, so this shows that $|h_0^{-1}(i)| = |F|$, as desired. (2) If $x \neq 0$, then for any choice of $b$, there is unique choice of $a$, namely $a = (i - b)x^{-1}$.

**¶17. Example:** (Random Graphs) Fix $0 \leq p \leq 1$ and $n \geq 2$. Consider the probability space where $\Omega = \{0, 1\}^m$, $m = \binom{n}{2}$, $\Sigma = 2^\Omega$ and for $(b_1, \ldots, b_m) \in \Omega$, $\Pr(b_1, \ldots, b_m) = p^k (1 - p)^{m-k}$ where $k$ is the number of 1's in $(b_1, \ldots, b_m)$. Once checks that Pr as defined is a probability function. Let $K_n$ be the complete bigraph on $n$ vertices whose edges are labelled with the integers $1, \ldots, m$. Consider the "random graph" which is just the function

$$G_{n,p} : \Omega \to \{\text{subgraphs of } K_n\} \tag{12}$$

where $G_{n,p}(b_1, \ldots, b_m)$ is the subgraph of $K_n$ such that the $i$th edge is in $G_{n,p}(b_1, \ldots, b_m)$ iff $b_i = 1$. Let us be concrete: let $n = 4$, $m = \binom{4}{2} = 6$, $p = 1/3$ and $G_4$ be the graph in the margin. Then $\Pr(G_4) = (1/3)^4 (2/3)^2 = 4/3^6$. Of course, $4/3^6$ is the probability of any subgraph of $K_4$ with 4 edges.



$G_4$

**¶18. Random Statistics.** Random variables often arise as follows. A function $C : D \to \mathbb{R}$ is called a **statistic** of $D$ where $D$ is some set of objects. If $g : \Omega \to D$ is a random object, we obtain the random variable $C_g : \Omega \to \mathbb{R}$ where

$$C_g(\omega) = C(g(\omega)).$$

Call $C_g$ a **random statistic** of $g$.

For example, let $g = G_{n,p}$ be the random graph in equation (12). and let $C$ count the number of Hamiltonian cycles in a bigraph. Then the random variable

$$C_g : \Omega \to \mathbb{R} \tag{13}$$

is defined so that $C_g(\omega)$ is the number of Hamiltonian cycles in $g(\omega)$.

**¶19. Independent Ensembles.**   A collection $K$ of random $D$-objects with a common underlying probability space is called an **ensemble** of $D$-objects. We extend the concepts of independent events to "independent ensembles".

A finite ensemble $\{X_1, X_2, \ldots, X_n\}$ of $n$ r.v.'s is **independent** if for all $m$-subset $\{i_1, i_2, \ldots, i_m\} \subseteq \{1, 2, \ldots, n\}$ and all $c_1, \ldots, c_m \in \mathbb{R}$, the events $\{X_{i_1} \leq c_1\}, \ldots, \{X_{i_m} \leq c_m\}$ are independent. Note that this is well-defined because the r.v.'s share a common probability space. An ensemble (finite or infinite) is $k$-**wise independent** if for all $m \leq k$, each $m$-subset of the ensemble is independent. As usual, **pairwise independent** means $2-wise$ independent.

Let $K$ be an ensemble of $D$-objects where $D$ is a finite set. We say $K$ is $k$-**wise independent** if for any $a_1, \ldots, a_k \in D$, and any $k$-subset $\{f_1, \ldots, f_k\} \subseteq K$, the events $\{f_1 = a_1\}, \ldots, \mathbf{f_k} = \mathbf{a_k}$ are independent.

**¶20. Example:**   (Finite Fields, contd) Recall the finite field space $\Omega = F^2$ above. Consider the ensemble

$$K = \{\mathbf{h}_x : x \in F\} \tag{14}$$

where $\mathbf{h}_x(\langle a, b \rangle) = ax + b$ as before. We had seen that the elements of $K$ are uniformly random. Let us show that the ensemble $K$ is pairwise independent: Fix $x, y, i, j \in F$ and let $n = |F|$. Suppose $x \neq y$ and $\mathbf{h}_x = i$ and $\mathbf{h}_y = j$. This means

$$\begin{pmatrix} x & 1 \\ y & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} i \\ j \end{pmatrix}.$$

The $2 \times 2$ matrix is invertible and hence $(a, b)$ has a unique solution. Hence

$$\Pr\{\mathbf{h}_x = i, \mathbf{h}_y = j\} = 1/n^2 = \Pr\{\mathbf{h}_x = i\}\Pr\{\mathbf{h}_y = j\},$$

as desired. Algorithmically, constructions of $k$-wise independent variables over a small sample space (e.g., $|\Omega| = p^2$ here) is important because it allows us to make certain probabilistic constructions effective.

_____EXERCISES

**Exercise 5.1:** Compute the probability of the event $\{C_g = 0\}$ where $C_g$ is given by (13). Do this for $n = 2, 3, 4$.                                                                                     ◇

**Exercise 5.2:** Let $\Omega = F$ be a finite field, and for $x \in F$, consider the random function $\mathbf{h}_x : \Omega \to F$ where $\mathbf{h}_x(a) = ax$. Consider the ensemble $K = \{\mathbf{h}_x : x \in F \setminus \{0\}\}$. Show that each $\mathbf{h}_x \in K$ is a uniformly random element of $F$, but $K$ is not pairwise independent.                                                                                     ◇

**Exercise 5.3:** Let $f_1, \ldots, f_n$ be real functions $f_i : \mathbb{R} \to \mathbb{R}$ and $\{X_1, \ldots, X_n\}$ is an independent ensemble of r.v.'s. If $f_i(X_i)$ are also r.v.'s then $\{f_1(X_1), \ldots, f_n(X_n)\}$ is also an independent ensemble.     ◇

**Exercise 5.4:** Consider the following (silly) randomized process, which is a sequence of probabilistic steps. At each step, we roll a dice that has one of six possible outcomes: $1, 2, 3, 4, 5, 6$. In the $i$-th step, if the outcome is less than $i$, we stop. Otherwise, we go to the next step. The first step is $i = 1$. For instance, we never stop after first step, and surely stop by the 7-th step. Let $T$ be the random variable corresponding to the number of steps.
(a) Set up the sample space, the event space, and the probability function for $T$.
(b) Compute the expected value of $T$.     ◇

**Exercise 5.5:** Let $K$ be an ensemble of random elements in the finite field $F$ (14).
(a) Show that $K$ is not 3-wise independent.
(b) Generalize the example to construct a collection of $k$-wise independent random functions.     ◇

**Exercise 5.6:** Let $W(n, x)$ (where $n \in \mathbb{N}$ and $x \in \mathbb{Z}_n$) be a "witness" predicate for compositeness: if $n$ is composite, then $W(n, x) = 1$ for at least $n/2$ choices of $x$; if $n$ is prime, then $W(n, x) = 0$ for all $x$. Let $W(n)$ be the random variable whose value is determined by a random choice of $x$. Let $W_t(n)$ be the random variable whose value is obtained as follows: randomly choose $n$ values $x_1, \ldots, x_n \in \mathbb{Z}_n$ and compute each $W(n, x_i)$. If any $W(n, x_i) = 1$ then $W_t(n) = 1$ but otherwise $W_t(n) = 0$.
(a) If $n$ is composite, what is the probability that $W_t(n) = 1$?
(b) Now we compute $W_t(n)$ using somewhat less randomness: first assume $t$ is prime and larger than $n$. only randomly choose two values $a, b \in \mathbb{Z}_t$. Then we define $y_i = a \cdot i + b \pmod{t}$. We evaluate $W_t(n)$ as before, except that we use $y_0, \ldots, y_{t-1} \pmod{n}$ instead of the $x_i$'s. Lower bound the probability that $W_t(n) = 1$ in this new setting.     ◇

_____END EXERCISES

# §6. Expectation and Variance

    Two important numbers are associated with a random variable: its "average value" and its "variance" (expected deviation from the average value). We had already defined the average value, $\mathrm{E}[X]$. Let us look at some of its properties.

**¶21. Two Remarkable Properties of Expectation.** We look at two elementary properties of expectation that often yield surprising consequences. The first is called **linearity of expectation**. This means that for all r.v.'s $X, Y$ and $\alpha, \beta \in \mathbb{R}$:

$$\mathrm{E}[\alpha X + \beta Y] = \alpha \mathrm{E}[X] + \beta \mathrm{E}[Y].$$

The proof for the discrete case is immediate: $\mathrm{E}[\alpha X + \beta Y] = \sum_\omega (\alpha X(\omega) + \beta Y(\omega)) \Pr(\omega) = \alpha \mathrm{E}[X] + \beta \mathrm{E}[Y]$. The remarkable fact is that $X$ and $Y$ are completely arbitrary – in particular,

there are no independence assumptions. In applications, we often decompose a r.v. $X$ into *some* linear combination of r.v.s $X_1, X_2, \ldots, X_m$. If we can compute the expectations of each $X_i$, then by linearity of expectation, we obtain the expectation of $X$ itself. Thus, $X$ may be the running time of a $n$-step algorithm and $X_i$ is the expected time for the $i$th step.

The second remarkable property is that, from expectations, we can infer the existence of objects with certain properties.

LEMMA 2. *Let $X$ be a discrete r. v. with finite expectation $\mu$. If $\Omega$ is finite, then:*

**(i)** *There exists $\omega_0, \omega_1 \in \Omega$ such that*

$$X(\omega_0) \leq \mu \leq X(\omega_1). \tag{15}$$

**(ii)** *If $X$ is non-negative and $c > 0$ then*

$$\Pr\{X < c\mu\} > 1/c. \tag{16}$$

*In particular, if $\Pr\{\cdot\}$ is uniform and $\Omega$ finite, then more than half of the sample points $\omega \in \Omega$ satisfy $X(\omega) < 2\mu$.*

*Proof.* Since $X$ is discrete, let

$$\mu = \mathrm{E}[X] = \sum_{i=1}^{\infty} a_i \Pr\{X = a_i\}.$$

(i) If there are arbitrarily negative $a_i$'s then clearly $\omega_0$ exists; otherwise choose $\omega_0$ so that $X(\omega_0) = \inf\{X(\omega) : \omega \in \Omega\}$. Likewise if there are arbitrarily large $a_i$'s then $\omega_1$ exists, and otherwise choose $\omega_1$ so that $X(\omega_1) = \sup\{X(\omega) : \omega \in \Omega\}$. In every case, we have chosen $\omega_0$ and $\omega_1$ so that the following inequality confirms our lemma:

$$X(\omega_0) = X(\omega_0) \sum_{\omega \in \Omega} \Pr(\omega) \leq \sum_{\omega \in \Omega} \Pr(\omega) X(\omega) \leq X(\omega_1) \sum_{\omega \in \Omega} \Pr(\omega) = X(\omega_1).$$

(ii) This is just Markov's inequality (below). We have $c\mu \cdot \Pr\{X \geq c\mu\} \leq \mathrm{E}[X] = \mu$, so $\Pr\{X \geq c\mu\} \leq \frac{1}{c}$.      **Q.E.D.**

To apply this lemma, we set up a random $D$ object,

$$g : \Omega \to D$$

and are interested in a certain statistic $C : D \to \mathbb{R}$. Define the random statistic $C_g : \Omega \to \mathbb{R}$ as in (13). Then there exists $\omega_0$ such that

$$C_g(\omega_0) \leq \mathrm{E}[C_g]$$

This means that the object $g(\omega_0) \in D$ has the property $C(g(\omega_0)) \leq \mathrm{E}[C_g]$.

Linearity of expectation amounts to saying that summing r.v.'s is commutative with taking expectation. What about products of r.v.'s? If $X, Y$ are independent then

$$\mathrm{E}[XY] = \mathrm{E}[X]\mathrm{E}[Y]. \tag{17}$$

The requirement that $X, Y$ be independent is necessary. As noted earlier, all multiplicative properties of probability depends from some form of independence.

The $j$th **moment** of $X$ is $\mathtt{E}[X^j]$. If $\mathtt{E}[X]$ is finite, then we define the **variance** of $X$ to be

$$\mathtt{Var}(X) := \mathtt{E}[(X - \mathtt{E}[X])^2].$$

Note that $X - \mathtt{E}[X]$ is the deviation of $X$ from its mean. It is easy to see that

$$\mathtt{Var}(X) = \mathtt{E}[X^2] - \mathtt{E}[X]^2.$$

The positive square-root of $\mathtt{Var}(X)$ is called its **standard deviation** and denoted $\sigma(X)$ (so $\mathtt{Var}(X)$ is also written $\sigma^2(X)$). If $X, Y$ are independent, then summing r.v.'s also commutes with taking variances. More generally:

LEMMA 3. *Let $X_i$ ($i = 1, \ldots, n$) be pairwise independent random variables with finite variances. Then*

$$\mathtt{Var}(\sum_i^n X_i) = \sum_{i=1}^n \mathtt{Var}(X_i).$$

This is a straightforward computation, using the fact that $\mathtt{E}[X_i X_j] = \mathtt{E}[X_i]\mathtt{E}[X_j]$ for $i \neq j$ since $X_i$ and $X_j$ are independent.

**¶22. Distribution and Density.** For any r.v. $X$, we define its **distribution function** to be $F_X : \mathbb{R} \to [0, 1]$ where

$$F_X(c) := \Pr\{X \leq c\}, \qquad c \in \mathbb{R}.$$

The importance of distribution functions stems from the fact that the basic properties of random variables can be studied from their distribution function alone.

Two r.v.'s $X, Y$ can be related as follows: we say $X$ **stochastically dominates** $Y$, written

$$X \succeq Y$$

if $F_X(c) \leq F_Y(c)$ for all $c$. This implies (Exercise) $\mathtt{E}[X] \geq \mathtt{E}[Y]$ if $X$ stochastically dominates $Y$. If $X \succeq Y$ and $Y \succeq X$ then we say they are **identically distributed**, denoted

$$X \sim Y.$$

A common probabilistic setting is an independent ensemble $K$ of r.v.'s that all share the same distribution. We then say $K$ is **independent and identically distributed** (abbrev. i.i.d) ensemble. For instance, when $X_i$ is the outcome of the $i$th toss of some fixed coin, then $K = \{X_i : i = 1, 2, \ldots\}$ is an i.i.d. ensemble.

In general, a distribution function[4] $F(x)$ is a monotone non-decreasing real function such that $F(-\infty) = 0$ and $F(+\infty) = 1$. Sometimes, a distribution function $F(x)$ is defined via a **density function** $f(u) \geq 0$, where

$$F(x) = \int_{-\infty}^x f(u)du.$$

In case $X$ is discrete, the density function $f_X(u)$ (of its distribution function $F_X$) is zero at all but countably many values of $u$. As defined above, a continuous r.v. $X$ is specified by its density function.

---

[4] Some authors calls the function $\Pr : \Sigma \to [0, 1]$ a "(probability) distribution" on the set $\Omega$. We avoid this terminology.

**¶23. Conditional Expectation.** This concept is useful for computing expectation. If $A$ is an event, define the **conditional expectation** $\mathbb{E}[X|A]$ of $X$ to be $\sum_{i \geq 1} a_i \Pr\{X = a_i|A\}$. In the discrete event space, we get

$$\mathbb{E}[X|A] = \frac{\sum_{\omega \in A} X(\omega)\Pr(\omega)}{\Pr(A)}.$$

If $B$ is the complement of $A$, then

$$\mathbb{E}[X] = \mathbb{E}[X|A]\Pr(A) + \mathbb{E}[X|B]\Pr(B).$$

More generally, if $Y$ is another r.v., we define a new r.v. $Z = \mathbb{E}[X|Y]$ where $Z(\omega) = \mathbb{E}[X|Y = Y(\omega)]$ for any $\omega \in \Omega$. We can compute the expectation of $X$ using the formula

$$\mathbb{E}[X] = \mathbb{E}[\mathbb{E}[X|Y]] \tag{18}$$
$$= \sum_{a \in \mathbb{R}} \mathbb{E}[X|Y = a]\Pr\{Y = a\}. \tag{19}$$

For example, let $X_i$'s be i.i.d., and $N$ be a non-negative integer r.v. independent of the $X_i$'s. What is the expected value of $\sum_{i=1}^{N} X_i$?

$$\mathbb{E}[\sum_{i=1}^{N} X_i] = \mathbb{E}[\mathbb{E}[\sum_{i=1}^{N} X_i|N]]$$
$$= \sum_{n \in \mathbb{N}} \mathbb{E}[\sum_{i=1}^{N} X_i|N = n]\Pr\{N = n\}$$
$$= \sum_{n \in \mathbb{N}} n\mathbb{E}[X_1]\Pr\{N = n\}$$
$$= \mathbb{E}[X_1]\mathbb{E}[N].$$

We can also use conditioning in computing variance, since $\mathbb{E}[X^2] = \mathbb{E}[\mathbb{E}[X^2|Y]]$.

_____Exercises

**Exercise 6.1:** Answer YES or NO to the following question. A correct answer is worth 5 points, but a wrong answer gets you $-3$ points. Of course, if you do not answer, you get 0 points. "In a True/False question, you get 5 points for correct answer, 0 points for not attempting the question and -3 points for an incorrect guess. Suppose have NO idea what the answer might be. Should you attempt to answer the question?".                    ◇

**Exercise 6.2:** You face a multiple-choice question with 4 possible choices. If you answer the question, you get 6 points if correct and -3 if wrong. If you do not attempt the question, you get -1 point. Should you attempt to answer the question if you have no clue as to what the question is about? You must justify your answer to receive any credit.                    ◇

_this_ is not a multiple choice question.

**Exercise 6.3:** You (as someone who is designing an examination) wants to assign points to a multiple choice question in which the student must pick one out of 5 possible choices. The student is not allowed to ignore the question. How do you assign points so that (a) if a student has no clue, then the expected score is $-1$ points and (b) if a student could eliminate one out of the 5 choices, the expected score is 0 points.                    ◇

**Exercise 6.4:** Compute the expected value of the r.v. $C_g$ in equation (13) for small values of $n$ ($n = 2, 3, 4, 5$).  ◇

**Exercise 6.5:** Simple dice game: you are charged $c$ dollars for rolling a dice, and if your roll has outcome $i$, you win $i$ dollars. What is the fair value of $c$? HINT: what is your expected win per roll?  ◇

**Exercise 6.6:** (a) Professor Vegas introduces a game of dice in class (strictly for "object lesson" of course). Anyone in class can play. To play the game, you pay $12 and roll a pair of dice. If the product of the rolled values on the dice is $n$, then Professor Vegas pays you $$n$. For instance, if you rolled the numbers 5 and 6 then you make a profit of $18 = 30 - 12$. Student Smart would not play, claiming: *the probability of losing money is more than the probability of winning money.*
(a) What is right and wrong with Student Smart's claim?
(b) Would you play this game? Justify.  ◇

**Exercise 6.7:** One day, Professor Vegas forgot to bring his pair of dice. He stills wants to play the game in the previous exercise. Professor Vegas decides to simulate the dice by tossing a fair coin 6 times. Interpreting heads as 1 and tails as zero, this gives 6 bits which can be viewed as two binary numbers $x = x_2 x_1 x_0$ and $y = y_2 y_1 y_0$. So $x$ and $y$ are between 0 and 7. If $x$ or $y$ is either 0 or 7 then the Professor returns your $12 (the game is off). Otherwise, this is like the dice game in (a). What is the expected profit of this game?  ◇

**Exercise 6.8:** In the previous question, we "simulate" rolling a dice by tossing three fair coins. Unfortunately, if the value of the tosses is 0 or 7, we call off the game. Now, we want to continue tossing coins until we get a value between 1 and 6.
(a) An obvious strategy is this: each time you get 0 or 7, you toss another three coins. This is repeated as many times as needed. What is the expected number of coin tosses to "simulate" a dice roll using this method?
(b) Modify the above strategy to simulate a dice roll with fewer coin tosses. You need to (i) justify that your new strategy simulates a fair dice and (ii) compute the expected number of coin tosses.
(c) Can you show what the optimum strategy is?  ◇

**Exercise 6.9:** In the dice game of the previous exercise, Student Smart decided to do another computation. He sets up a sample space

$$S = \{11, 12, \ldots, 16, 22, 23, \ldots, 26, 33, \ldots, 36, 44, 45, 46, 55, 56, 66\}.$$

So $|S| = 21$. Then he defines the r.v. $X$ where $X(ij) = i \times j$ and computes the expectation of $X$ where using $\Pr(ij) = 1/21$. What is wrong? Can you correct his mistake without changing his choice of sample space? What is the alternative sample space? In what sense is Smart's choice of $S$ is better?  ◇

**Exercise 6.10:** In this game, you begin with a pot of 0 dollars. Before each roll of the dice, you can decide to stop or to continue. If you decide to stop, you keep the pot. If you decide to play, and if you roll any value less than 6, that many dollars are added to the pot. But if you roll a 6, you lose the pot and the game ends. What is your optimal strategy? What if you must pay $Z$ dollars to play?  ◇

**Exercise 6.11:** Let $X, Y$ takes on finitely many values $a_1, \ldots, a_n$. (So they are discrete r.v.'s) Prove if $X \succeq Y$ then $\mathtt{E}[X] \geq \mathtt{E}[Y]$ and equality holds iff $X \sim Y$.                    ◇

**Exercise 6.12:** (a) Show that in any graph with $n$ vertices and $e$ edges, there exists a bipartite subgraph with $e/2$ edges. In addition, the bipartite subgraph have $\lfloor n \rfloor$ vertices on one side and $\lceil n \rceil$ of the other. Remark: depending on your approach, you may not be able to fulfill the additional requirement.
(b) Obtain the same result constructively (*i.e.*, give a randomized algorithm).                    ◇

**Exercise 6.13:** (Cauchy-Schwartz Inequality) Show that $\mathtt{E}[XY]^2 \leq \mathtt{E}[X^2]\mathtt{E}[Y^2]$ assuming $X, Y$ have finite variances.                    ◇

**Exercise 6.14:** (Law of Unconscious Statistician) If $X$ is a discrete r.v. with probability mass function $f_X(u)$, and $g$ is a real function then

$$\mathtt{E}[g(X)] = \sum_{u:f_X(u)>0} g(u)f_X(u).$$

◇

**Exercise 6.15:** If $X_1, X_2, \ldots$ are i.i.d. and $N \geq 0$ is an independent r.v. that is integer-valued then $\mathtt{E}[\sum_{i=1}^{N} X_i] = \mathtt{E}[N]\mathtt{E}[X_1]$ and $\mathtt{Var}(\sum_{i=1}^{N} X_i) = \mathtt{E}[N]\mathtt{Var}(X_1) + \mathtt{E}[X]^2\mathtt{Var}(N)$.                    ◇

**Exercise 6.16:** Suppose we have a fair game in which you can bet any dollar amount. If you bet \$x, and you win, you receive \$x; and otherwise you lose \$x.
(a) A well-known "gambling technique" is to begin by betting \$1. Each time you lose, you double the amount of the bet (to \$2, \$4, etc). You stop at the first time you win. What is wrong with this scenario?
(b) Suppose you have a limited amount of dollars, and you want to devise a strategy in which the *probability* of your winning is as big as possible. (We are not talking about your "expected win".) How would you achieve this?                    ◇

**Exercise 6.17:** [Amer. Math. Monthly] A set consisting of $n$ men and $n$ women are partitioned at random into $n$ disjoint pairs of people. Let $X$ be the number of male-female couples that result. What is the expected value and variance of $X$? HINT: let $X_i$ be the indicator variable for the event that the $i$th man is paired with a woman. To compute the variance, first compute $\mathtt{E}[X_i^2]$ and $\mathtt{E}[X_i X_j]$ for $i \neq j$.                    ◇

**Exercise 6.18:** [Mean and Variance of a geometric distribution] Let $X$ be the number of coin tosses needed until the first head appears. Assume the probability of coming up heads is $p$. Use conditional probability (18) to compute $\mathtt{E}[X]$ and $\mathtt{Var}(X)$. HINT: let $Y = 1$ if the first toss is a head, and $Y = 0$ else.                    ◇

_____ END EXERCISES

## §7. Analysis of QuickSort

In Chapter II, we gave an analysis of Quicksort by solving a recurrence. We now re-visit the analysis, using a more probabilistic language.

Assume the input is an array $A[1..n]$ holding $n$ numbers. Recall that Quicksort picks a random $r \in \{1, \dots, n\}$ and uses the value $A[r]$ to partition the numbers in $A[1..n]$ into those that are greater than $A[r]$ and those that are less than $A[r]$.

We first describe an elegant partition subroutine that requires no extra array storage beyond the input array $A$:

```
PARTITION(A, i, j, e):
Input: Array A[1..n] and 1 ≤ i < j ≤ n
Output: Index k ∈ {i,...,j} and a rearranged subarray A[i..j] satisfying
        A[i..k] ≤ e and A[k+1..j] > e.
    While (j − i ≥ 1)
        While (i < j and A[i] ≤ e) i++;
        While (i < j and A[j] ≥ e) j--;
        If (i < j)
            Swap A[i] ↔ A[j];   i++;   j--
    ▷ At this point, i = j or j + 1 = i
    If (j = i)
        If (A[i] ≤ e) Return (i)
        else Return (i − 1)
    else
        Return (j)
```

To sort the array $A[1..n]$, we invoke QUICKSORT$(A, 1, n)$, where QUICKSORT is the following recursive procedure:

```
QUICKSORT(A, i, j):
Input: Array A[1..n] and 1 ≤ i ≤ j ≤ n.
Output: The subarray A[i..j] is sorted in non-decreasing order.
    If i = j, Return
    Randomly pick a r ∈ {i,...,j}
    Swap A[r] ↔ A[i]
    k ←PARTITION(A, i, j, A[i])
    Swap A[i] ↔ A[k]
    If (i < k) QUICKSORT(A, i, k − 1)
    If (k < j) QUICKSORT(A, k + 1, j)
```

To simplify the analysis, we assume the input numbers are distinct. Suppose the set of numbers in $A[1..n]$ is
$$Z := \{z_1, \dots, z_n\}$$
where $z_1 < z_2 < \cdots < z_n$. For each $1 \le i < j \le n$, let

$$E_{ij} = \{z_i \text{ is compared to } z_j\}$$

denote the event that that $z_i$ and $z_j$ are compared. Let $X_{ij}$ be the indicator function for $E_{ij}$.

If $X$ is the random variable for the number of comparisons in quicksort, then we have

$$X = \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} X_{ij}.$$

The critical observation is:

Lemma 4. $\Pr(E_{ij}) = \frac{2}{j-i+1}$

From this lemma, we see that

$$\mathrm{E}[X] = \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} \mathrm{E}[X_{ij}] = \sum_{i=1}^{n-1} \sum_{j=i+1}^{n} \Pr(E_{ij}) < \sum_{i=1}^{n-1} 2H_n < 2n \lg n.$$

It remains to prove Lemma 4. We will prove a more general statement: let $A_{ij}$ be the event that the pivot $r$ on input $Z = \{1, \ldots, n\}$ satisfies the property $r < i$ or $r > j$. For all $1 \le i < j \le n$, let $E_{ij}^n$ denote the event that the comparison $z_i : z_j$ occurred on input $Z = \{1, \ldots, n\}$. When $i = 1$ and $j = n$, we let

$$E^n := E_{1,n}^n.$$

We claim: if $1 < i$ or $j < n$ then

$$\Pr(E_{ij}^n | A_{ij}) = \Pr(E^{j-i+1}). \tag{20}$$

Since $1 < i$ or $j < n$, we must have $n - j + i - 1 \ge 1$. We use induction on $n - j + i - 1$. The result is clearly true when $n - j + i - 1 = 1$. Otherwise, $n - j + i = 1 > 1$ and we see that

$$
\begin{aligned}
\Pr(E_{ij}^n | A_{ij}) \cdot \Pr(A_{ij}) &= \Pr(E_{ij}^n A_{ij}) \\
&= \Pr(E_{ij}^n | r < i) \Pr\{r < i\} + \Pr(E_{ij}^n | r > j) \Pr\{r > j\} \\
&= \Pr(E^{j-i+1}) \Pr\{r < i\} + \Pr(E^{j-i+1}) \Pr\{r > j\} \qquad \text{(by induction hypothesis)} \\
&= \Pr(E^{j-i+1}) \left[ \tfrac{i-1}{n} + \tfrac{n-j}{n} \right] \\
&= \Pr(E^{j-i+1}) \Pr(A_{ij}).
\end{aligned}
$$

But it is immediate that

$$\Pr(E^n) = \frac{2}{n}$$

and thus $\Pr(E_{ij}^n) = \Pr(E^{j-i+1}) = \frac{2}{j-i+1}$. This proves Lemma 4.

_____Exercises

**Exercise 7.1:** Let us estimate the size $C(n)$ of the sample space for Quicksort of $n$ numbers.
(a) Show that $C(n)$ satisfy the recurrence $C(n) = \sum_{i=1}^{n} C(i-1)C(n-i)$ for $n \ge 1$ and $C(0) = 1$.
(b) If $G(x) = \sum_{n=0}^{\infty} C(n)x^n$ is the generating function of the $C(n)$'s, show that $G(x) = (1 - \sqrt{1 - 4x})/2x$. HINT: What is the connection between $G(x)$ and $G(x)^2$?
(c) Use the Taylor expansion of $\sqrt{1 - 4x}/2x$. at $x = 0$ to show

$$C(n) = \frac{1}{n+1} \binom{2n}{n}. \tag{21}$$

Then use Stirling's approximation (Chapter II) to give tight bounds on $C(n)$.    ◇

_____End Exercises

## §8. Ensemble of Random Variables

We consider ensembles of random variables in two common situations.

(i) An ensemble of i.i.d. r.v.'s.

(ii) An ensemble $\{X_t : t \in T\}$ of r.v.'s where $T \subseteq \mathbb{R}$ is the index set. We think of $T$ as time and $X_t$ as describing the behavior of a stochastic phenomenon evolving over time. Such a family is called a **stochastic process**. Usually, we have either $T = \mathbb{R}$ (continuous time) or $T = \mathbb{N}$ (discrete time).

We state two results that lay claim to being the fundamental theorems of probability theory. Both relate to i.i.d. ensembles. Let $X_1, X_2, X_3, \dots$, be a countable i.i.d. ensemble of Bernoulli r.v.'s. Let

$$S_n := \sum_{i=1}^{n} X_i, \quad p := \Pr\{X_1 = 1\}, \quad \sigma := \sqrt{\mathtt{Var}(X_i)}.$$

It is intuitively clear that $S_n$ approaches $np$ as $n \to \infty$.

THEOREM 5 (Strong Law of Large Numbers). *For any $\varepsilon > 0$, with probability $1$, there are only finitely many sample points in the event*

$$\{|S_n - np| > \varepsilon\}$$

THEOREM 6 (Central Limit Theorem).

$$\Pr\left\{ \frac{S_n - np}{\sigma\sqrt{n}} \le t \right\} \to \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{t} e^{-x^2/2} dx$$

*as $n \to \infty$.*

Alternatively, the Central Limit Theorem says that the distribution of $\frac{S_n - np}{\sigma\sqrt{n}}$ tends to the standard normal distribution (see below) as $n \to \infty$.

¶24. **Some probability distributions.** The above theorems do not make any assumptions about the underlying distributions of the r.v.'s (therein lies their power). However, certain probability distributions are quite common and it is important to recognize them. Below we list some of them. In each case, we only need to describe the corresponding density functions $f(u)$. In the discrete case, it suffices to specify $f(u)$ at those elementary events $u$ where $f(u) > 0$.

- **Binomial distribution** $B(n, p)$, with parameters $n \ge 1$ and $0 < p < 1$:

$$f(i) := \binom{n}{i} p^i (1-p)^{n-i}, \qquad (i = 0, 1, \dots, n).$$

Sometimes $f(i)$ is also written $B_i(n, p)$ and corresponds to the probability of $i$ successes out of $n$ Bernoulli trials. In case $n = 1$, this is also called the Bernoulli distribution. If $X$ has such a distribution, then

$$\mathtt{E}[X] = np, \quad \mathtt{Var}(X) = npq$$

where $q = 1 - p$.

- **Geometric distribution** with parameter $p$, $0 < p < 1$:

$$f(i) := p(1-p)^{i-1} = pq^{i-1}, \qquad (i = 1, 2, \ldots).$$

Thus $f(i)$ may be interpreted as the probability of the first success occurring at the $i$th Bernoulli trial. If $X$ has such a distribution, then $\mathtt{E}[X] = 1/p$ and $\mathtt{Var}(X) = q/p^2$.

- **Poisson distribution** with parameter $\lambda > 0$:

$$f(i) := e^{-\lambda} \frac{\lambda^i}{i!}, \qquad (i = 0, 1, \ldots).$$

We may view $f(i)$ as the limiting case of $B_i(n, p)$ where $n \to \infty$ and $np = \lambda$. If $X$ has such a distribution, then $\mathtt{E}[X] = \mathtt{Var}(X) = \lambda$.

- **Uniform distribution** over the real interval $[a, b]$:

$$f(u) := \begin{cases} \frac{1}{b-a} & a < u < b \\ 0 & \text{else.} \end{cases}$$

- **Exponential distribution** with parameter $\lambda > 0$:

$$f(u) = \begin{cases} \lambda e^{-\lambda u} & u \geq 0 \\ 0 & \text{else.} \end{cases}$$

- **Normal distribution** with mean $\mu$ and variance $\sigma^2$:

$$f(u) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[ -\frac{1}{2}\left( \frac{u - \mu}{\sigma} \right)^2 \right].$$

In case $\mu = 0$ and $\sigma^2 = 1$, we call this the unit normal distribution.

—————————————————————————————EXERCISES

**Exercise 8.1:** Verify the values of $\mathtt{E}[X]$ and $\mathtt{Var}(X)$ asserted for the various distributions of $X$.

$\diamondsuit$

**Exercise 8.2:** Show that the density functions $f(u)$ above truly define distribution functions: $f(u) \geq 0$ and $\int_{-\infty}^{\infty} f(u)du = 1$. Determine the distribution function in each case. $\diamondsuit$

—————————————————————————————END EXERCISES

## §9. Estimates and Inequalities

A fundamental skill in probabilistic analysis is estimating probabilities because they are often too intricate to determine exactly. We list some useful inequalities and estimation techniques.

**¶25. Approximating the binomial coefficients.** Recall Stirling's approximation in Lecture II.2. Using such bounds, we can show [13] that for $0 < p < 1$ and $q = 1 - p$,

$$G(p,n)e^{-\frac{1}{12pn} - \frac{1}{12qn}} < \binom{n}{pn} < G(p,n) \tag{22}$$

where

$$G(p,n) = \frac{1}{\sqrt{2\pi pqn}} p^{-pn} q^{-qn}.$$

**¶26. Tail of the binomial distribution.** The "tail" of the distribution $B(n,p)$ is the following sum

$$\sum_{i=\lambda n}^{n} \binom{n}{i} p^i q^{n-i}.$$

It is easy to see the following inequality:

$$\sum_{i=\lambda n}^{n} \binom{n}{i} p^i q^{n-i} \le \binom{n}{\lambda n} p^{\lambda n}.$$

To see this, note that LHS is the probability of the event $A = \{$There are at least $\lambda n$ successes in $n$ coin tosses$\}$. For any choice $x$ of $\lambda n$ out of $n$ coin tosses, let $B_x$ be the event that the chosen coin tosses are successes. Then the LHS is the sum of the probability of $B_x$, over all $x$. Clearly $A = \cup_x B_x$. But the RHS may be an over count because the events $B_x$ need not be disjoint. We have the following upper bound [5]:

$$\sum_{i=\lambda n}^{n} \binom{n}{i} p^i q^{n-i} < \frac{\lambda q}{\lambda - p} \binom{n}{\lambda n} p^{\lambda n} q^{\mu n}$$

where $\lambda > p$ and $q = 1 - p$. This specializes to

$$\sum_{i=\lambda n}^{n} \binom{n}{i} < \frac{\lambda}{2\lambda - 1} \binom{n}{\lambda n}$$

where $\lambda > p = q = 1/2$.

**¶27. Markov Inequality.** Let $X$ be a non-negative random variable. We have the trivial bound

$$\Pr\{X \ge 1\} \le \mathtt{E}[X]. \tag{23}$$

For any real constant $c > 0$,

$$
\begin{aligned}
\Pr\{X \ge c\} &= \Pr\{X/c \ge 1\} \\
&\le \mathtt{E}[X/c] &&\text{(by trivial bound (23))} \\
&= \mathtt{E}[X]/c.
\end{aligned}
$$

This proves the **Markov inequality**,

$$\Pr\{X \ge c\} \le \frac{\mathtt{E}[X]}{c}. \tag{24}$$

Another proof uses the **Heaviside function** $H(x)$ that is the 0-1 function given by $H(x) = 1$ if and only if $x > 0$. We have the trivial inequality $H(X - c) \le \frac{X}{c}$ $(c > 0)$. Taking expectations on both sides yields the Markov inequality since $\mathtt{E}[H(X - c)] = \Pr\{X \ge c\}$.

For instance, we infer that the probability that $X$ is twice the expected value is at most $1/2$. Observe that the Markov inequality is trivial unless $\mathtt{E}[X]$ is finite and we choose $c > \mathtt{E}[X]$.

**¶28. Chebyshev Inequality.** It is also called the Chebyshev-Bienaymé inequality since it originally appeared in a paper of Bienaymé in 1853 [8, p. 73]. With any real $c > 0$,

$$\Pr\{|X| \geq c\} = \Pr\{X^2 \geq c^2\} \leq \frac{\mathtt{E}[X^2]}{c^2} \tag{25}$$

by an application of Markov inequality. In contrast to Markov's inequality, Chebyshev does not need to assume that $X$ is non-negative. The price is that we now need to take the expectation of $X^2$, making this bound somewhat "less elementary". Another form of this inequality (derived in exactly the same way) is

$$\Pr\{|X - \mathtt{E}[X]| \geq c\} = \Pr\{(X - \mathtt{E}[X])^2 \geq c^2\} \leq \frac{\mathtt{Var}(X)}{c^2}. \tag{26}$$

Sometimes $\Pr\{|X - \mathtt{E}[X]| \geq c\}$ is called the **tail probability** of $X$. By a trivial transformation of parameters, equation (26) can also written as

$$\Pr\{|X - \mathtt{E}[X]| \geq c\sqrt{\mathtt{Var}(X)}\} \leq \frac{1}{c^2}. \tag{27}$$

This form is useful in statistics because it bounds the probability of $X$ deviating from its mean by some fraction of the standard deviation, $\sqrt{\mathtt{Var}(X)}$.

Let us give an application of Chebyshev's inequality:

LEMMA 7. *Let $X$ be a r.v. with mean $\mathtt{E}[X] = \mu \geq 0$.*
*(a) Then*

$$\Pr\{X = 0\} \leq \frac{\mathtt{Var}(X)}{\mu^2}.$$

*(b) Suppose $X = \sum_{i=1}^{n} X_i$ where the $X_i$'s are pairwise independent Bernoulli r.v.s with $\mathtt{E}[X_i] = p$ (and $q = 1 - p$) then*

$$\Pr\{X = 0\} \leq \frac{q}{np}.$$

*Proof.* (a) Since $\{X = 0\} \subseteq \{|X - \mu| \geq \mu\}$, we have

$$\Pr\{X = 0\} \leq \Pr\{|X - \mu| \geq \mu\} \leq \frac{\mathtt{Var}X}{\mu^2}$$

by Chebyshev.
(b) It is easy to check that $\mathtt{Var}(X_i) = pq$. Since the $X_i$'s are independent, we have $\mathtt{Var}(X) = npq$. Also $\mathtt{E}[X] = \mu = np$. Plugging into the formula in (a) yields the claimed bound on $\Pr\{X = 0\}$.      **Q.E.D.**

Part (b) is useful in reducing the error probability in a certain class of randomized algorithms called $RP$-algorithms. The outcome of an $RP$-algorithm $A$ may be regarded as a Bernoulli r.v. $X_i$ which has value 1 or 0. If $X_i = 1$, then the algorithm has no error. If $X_i = 0$, then the probability of error is at most $p$ ($0 \leq p < 1$). We can reduce the error probability in $RP$-algorithms by repeating its computation $n$ times and output 0 iff each of the $n$ repeated computations output 0. Then part (b) bounds the error probability of the iterated computation. We will see several such algorithms later (e.g., primality testing in §XIX.2).

**¶29. Jensen's Inequality.** Let $f(x)$ be a real function. By definition, $f$ is **convex** means that for all $n$,

$$f\left(\sum_{i=1}^{n} p_i x_i\right) \leq \sum_{i=1}^{n} p_i f(x_i) \tag{28}$$

where $\sum_{i=1}^{n} p_i = 1$ and $p_i \geq 0$. This is equivalent to the case where $n = 2$. If $X$ and $f(X)$ are random variables then

$$f(\text{E}[X]) \leq \text{E}[f(X)].$$

Let us prove this for the case when $X$ has takes on finitely many values $x_i$ with probability $p_i$. Then $\text{E}[X] = \sum_i p_i x_i$ and

$$f(\text{E}[X]) = f(\sum_i p_i x_i) \leq \sum_i p_i f(x_i) = \text{E}[f(X)].$$

For instance, if $r \geq 1$ then $\text{E}[|X|^r] \geq (\text{E}[|X|])^r$.

_____EXERCISES

**Exercise 9.1:** Verify the equation (22).                                          ◇

**Exercise 9.2:** Describe the class of non-negative random variables for which Markov's inequality is tight.                                          ◇

**Exercise 9.3:** Chebyshev's inequality is the best possible. In particular, show an $X$ such that $\Pr\{|X - \text{E}[X]| > e\} = \text{Var}(X)/e^2$.                                          ◇

## §10. Chernoff Bounds

Suppose we wish an upper bound on the probability $\Pr\{X \geq c\}$ where $X$ is an arbitrary r.v.. To apply Markov's inequality, we need to convert $X$ to a non-negative r.v. One way is to use the r.v. $X^2$, as in the proof of Chebyshev's inequality. The technique of Chernoff converts $X$ to the Markov situation by using

$$\Pr\{X \geq c\} = \Pr\{e^X \geq e^c\}.$$

Since $e^X$ is a non-negative r.v., we conclude from Markov's inequality (24) that

$$\Pr\{X \geq c\} \leq e^{-c}\text{E}[e^X]. \tag{29}$$

We can exploit this trick further: for any positive number $t > 0$, we have $\Pr\{X \geq c\} = \Pr\{tX \geq tc\}$, and proceeding as before, we obtain

$$\begin{aligned}\Pr\{X \geq c\} &\leq& e^{-ct}\text{E}[e^{tX}] \\ &=& \text{E}[e^{t(X-c)}].\end{aligned}$$

We then choose $t$ to minimize the right-hand side of this inequality:

LEMMA 8 (Chernoff Bound). *For any r.v. $X$ and real $c$,*

$$\Pr\{X \geq c\} \leq m(c). \tag{30}$$

*where*

$$m(c) = m_X(c) := \inf_{t>0} \text{E}[e^{t(X-c)}]. \tag{31}$$

The name "Chernoff bound" [3] is applied to any inequality derived from (30). We now derive such Chernoff bounds under various scenarios.

Let $X_1, \ldots, X_n$ be independent and

$$S = X_1 + \cdots + X_n.$$

It is easily verified that then $e^{tX_1}, \ldots, e^{tX_n}$ (for any constant $t$) are also independent. Then equation (17) implies

$$\mathrm{E}[e^{tS}] = \mathrm{E}[\prod_{i=1}^{n} e^{tX_i}] = \prod_{i=1}^{n} \mathrm{E}[e^{tX_i}].$$

(A) Suppose that, in addition, the $X_1, \ldots, X_n$ are i.i.d., and $m(c)$ is defined as in (31). This shows

$$
\begin{aligned}
\Pr\{S \geq nc\} &\leq e^{-nct}\mathrm{E}[e^{tS}], \quad (t > 0)\\
&\leq [m(c)]^n.
\end{aligned}
$$

This is a generalization of (30).

(B) Assume $S$ has the distribution $B(n,p)$. It is not hard to compute that

$$m(c) = \left(\frac{p}{c}\right)^c \left(\frac{1-p}{1-c}\right)^{1-c}. \tag{32}$$

Then for any $0 < \varepsilon < 1$:

$$\Pr\{S \geq (1-\varepsilon)np\} \leq \left(\frac{1}{1-\varepsilon}\right)^{(1-\varepsilon)np} \left(\frac{1-p}{1-(1-\varepsilon)p}\right)^{n-(1+\varepsilon)np}.$$

We still need to make this bound more convenient for application:

$$
\begin{aligned}
\Pr\{S \geq (1+\varepsilon)np\} &\leq \exp(-\varepsilon^2 np/3) \tag{33}\\
\Pr\{S \leq (1-\varepsilon)np\} &\leq \exp(-\varepsilon^2 np/2) \tag{34}
\end{aligned}
$$

$$\tag{35}$$

Need the $\leq$ and $\geq$ version of Chernoff bound...          INCOMPLETE

(C) Now suppose the $X_i$'s are[5] independent Bernoulli variables with $\Pr\{X_i = 1\} = p_i$ ($0 \leq p_i \leq 1$) and $\Pr\{X_i = 0\} = 1 - p_i$ for each $i$. Then

$$\mathrm{E}[X_i] = p_i, \qquad \mu := \mathrm{E}[S] = \sum_{i=1}^{n} p_i.$$

Fix any $\delta > 0$. Then

$$
\begin{aligned}
\Pr\{S \geq (1+\delta)\mu\} &\leq m((1+\delta)\mu)\\
&= \inf_{t>0} \mathrm{E}[e^{t(X-(1+\delta)\mu)}]\\
&= \inf_{t>0} \frac{\mathrm{E}[e^{tX}]}{e^{(1+\delta)\mu}}.
\end{aligned}
$$

---

[5] This can be regarded as a generalization of Bernoulli trials where all the $p_i$'s are the same. This generalization is sometimes known as Poisson trials.

We leave as an exercise to optimize the $t$ to obtain

$$\Pr\{S \geq (1+\delta)\mu\} \leq \left(\frac{e^{\delta}}{(1+\delta)^{1+\delta}}\right)^{\mu}. \tag{36}$$

We can also use this technique to bound the probability that $S$ is smaller than $(1-\delta)\mu$: here we obtain

$$\Pr\{S \leq (1+\delta)\mu\} \leq e^{-\mu\delta^2/2}. \tag{37}$$

**¶30. Estimating a Probability and Hoeffding Bound.**   Consider the natural problem of estimating $p$ $(0 < p < 1)$ where $p$ is the probability that a given coin will show up heads in a toss. The obvious solution is to choose some reasonably large $n$, toss this coin $n$ times, and estimate $p$ by the ratio $h/n$ where $h$ is the number of times we see heads in the $n$ coin tosses.

This problem is still not well-defined since we have no constraints on $n$. So assume our goal is to satisfy the bound

$$\Pr\{|p - (h/n)| > \delta\} \leq \varepsilon \tag{38}$$

where $\delta$ is the **precision parameter** and $\varepsilon$ is a bound on the **error probability**. Given $\delta$ and $\varepsilon$, $(0 < \delta, \varepsilon < 1)$, we now have a well-defined problem. This problem seems to be solved by the Chernoff bound (B) in (33) where $S = X_1 + \cdots X_n$ is now interpreted to be $h$. Then

$$\{|p - (h/n)| > \delta\} = \{|np - h| > n\delta\} = \{|np - S| > n\delta\}$$

If we substitute $\delta$ with $p\varepsilon$, then we obtain

$$
\begin{aligned}
\Pr\{|p - (h/n)| > \delta\} &= \Pr\{|np - S| > n\delta\} \\
&= \Pr\{|S - np| > np\varepsilon\} \\
&\leq 2\exp(-
\end{aligned}
$$

The problem is that the $p$ that we are estimating appears on the right hand side. Instead, we need the following Hoeffding bound:

$$
\begin{aligned}
\Pr\{S > np + \delta\} &\leq \exp(-n\delta^2/2) &(39) \\
\Pr\{S < np - \delta\} &\leq \exp(-n\delta^2/2) &(40) \\
\Pr\{|S < np| > \delta\} &\leq 2\exp(-n\delta^2/2) &(41)
\end{aligned}
$$

Comparing the usual Chernoff bounds with the Hoeffding bound, we see that the former bound the relative error in the estimate while the latter concerns absolute error.

For a survey of Chernoff Bounds, see T. Hagerub and C. Rüb, "A guided tour of Chernoff Bounds", **Information Processing Letters** 33(1990)305–308.

―――――――――――――――――――――――――――――――Exercises

**Exercise 10.1:** Verify the equation (32).        ◇

**Exercise 10.2:** Obtain an upper bound on $\Pr\{X \leq c\}$ by using Chernoff's technique. HINT: $\Pr\{X \leq c\} = \Pr\{tX \geq tc\}$ where $t < 0$.        ◇

**Exercise 10.3:** Show the following:
    i) Bonferroni's inequality,

$$\Pr(AB) \geq \Pr(A) + \Pr(B) - 1.$$

    ii) Boole's inequality,

$$\Pr(\cup_{i=1}^{n} A_i) \leq \sum_{i=1}^{n} \Pr(A_i).$$

(This is trivial, and usually used without acknowledgment.) iii) For all real $x$, $e^{-x} \geq 1 - x$ with equality only if $x = 0$.
iv) $1 + x < e^x < 1 + x + x^2$ which is valid for $|x| < 1$.

$\Diamond$

**Exercise 10.4:** Kolmogorov's inequality: let $X_1, \ldots, X_n$ be mutually independent with expectation $\mathtt{E}[X_i] = m_i$ and variance $\mathtt{Var}(X_i) = v_i$. Let $S_i = X_1 + \cdots + X_i$, $M_i = \mathtt{E}[S_i] = m_1 + \cdots + m_i$ and $V_i = \mathtt{Var}(S_i) = v_1 + \cdots + v_i$. Then for any $t > 0$, the probability that the $n$ inequalities

$$|S_i - M_i| < tV_n, \qquad i = 1, \ldots, n,$$

holds simultaneously is at least $1 - t^{-2}$.      $\Diamond$

**Exercise 10.5:** (Motwani-Raghavan) (a) Prove the Chernoff bound (36). HINT: Show that $\mathtt{E}[e^{tX_i}] = 1 + p_i(e^t - 1)$ and $1 + x < e^x$ $(x = p_i(e^t - 1))$ implies $\mathtt{E}[e^{tS}] = e^{(e^t-1)\mu}$. Choose $t = \ln(1 + \delta)$ to optimize.
(b) Prove the Chernoff bound (36). HINT: Reduce to the previous situation using $\Pr\{X \leq c\} = \Pr\{-X \geq -c\}$. Also, $(1 - \delta)^{1-\delta} > e^{-\delta+\delta^2/2}$.      $\Diamond$

**Exercise 10.6:** We want to process a sequence of **requests** on a single (initially empty) list. Each request is either an insertion of a key or the lookup on a key. The probability that any request is an insertion is $p$, $0 < p < 1$. The cost of an insertion is 1 and the cost of a lookup is $m$ if the current list has $m$ keys. After an insertion, the current list contains one more key.
(a) Compute the expected cost to process a sequence of $n$ requests.
(b) What is the *approximate* expected cost to process the $n$ requests if we use a binary search tree instead? Assume that the cost of insertion, as well as of lookup, is $\log_2(1+m)$ where $m$ is the number of keys in the current tree. NOTE: If $L$ is a random variable (say, representing the length of the current list), assume that $\mathtt{E}[\log_2 L] \approx \log_2 \mathtt{E}[L]$, (*i.e.*, the expected value of the log is approximately the log of the expected value).
(c) Let $p$ be fixed, $n$ varying. Describe a rule for choosing between the two data structures. Assuming $n \gg 1 \gg p$, give some rough estimates (assume $\ln(n!)$ is approximately $n \ln n$ for instance).
(d) Justify the approximation $\mathtt{E}[\log_2 L] \approx \log_2 \mathtt{E}[L]$ as reasonable.      $\Diamond$

_____ END EXERCISES

## §11. **Generating Functions**

> In this section, we assume that our r.v.'s are discrete with range $\mathbb{N} = \{0, 1, 2, \ldots\}$.

This powerful tool of probabilistic analysis was introduced by Euler (1707-1783). If $a_0, a_1, \ldots$, is a denumerable sequence of numbers, then its **(ordinary) generating function** is the power series

$$G(t) := a_0 + a_1 t + a_2 t^2 + \cdots = \sum_{i=0}^{\infty} a_i t^i.$$

If $a_i = \Pr\{X = i\}$ for $i \geq 0$, we also call $G(t) = G_X(t)$ the **generating function of** $X$. We will treat $G(t)$ purely formally, although under certain circumstances, we can view it as defining a real (or complex) function of $t$. For instance, if $G(t)$ is a generating function of a r.v. $X$ then $\sum_{i \geq 0} a_i = 1$ and the power series converges for all $|t| \leq 1$. The power of generating functions comes from the fact that we have a compact packaging of a potentially infinite series, facilitating otherwise messy manipulations. Differentiating (formally),

$$\begin{aligned} G'(t) &= \sum_{i=1}^{\infty} i a_i t^{i-1}, \\ G''(t) &= \sum_{i=2}^{\infty} i(i-1) a_i t^{i-2}. \end{aligned}$$

If $G(t)$ is the generating function of $X$, then

$$G'(1) = \mathtt{E}[X], \qquad G''(1) = \mathtt{E}[X^2] - \mathtt{E}[X].$$

It is easy to see that if $G_1(t) = \sum_{i \geq 0} a_i t^i$ and $G_2(t) = \sum_{i \geq 0} b_i t^i$ are the generating functions of independent r.v.'s $X$ and $Y$ then

$$G_1(t) G_2(t) = \sum_{i \geq 0} t^i \sum_{j=0}^{i} a_j b_{i-j} = \sum_{i \geq 0} t^i c_i$$

where $c_i = \Pr\{X + Y = i\}$. Thus we have: the product of the generating functions of two independent random variables $X$ and $Y$ is equal to the generating function of their sum $X + Y$. This can be generalized to any finite number of independent random variables. In particular, if $X_1, \ldots, X_n$ are $n$ independent coin tosses (running example (E1)), then the generating function of $X_i$ is $G_i(t) = q + pt$ where $q := 1 - p$. So the generating function of the r.v. $S_n := X_1 + X_2 + \cdots + X_n$ is

$$(q + pt)^n = \sum_{i=0}^{n} \binom{n}{i} p^i q^{n-i} t^i.$$

Thus, $\Pr\{S_n = i\} = \binom{n}{i} p^i q^{n-i}$ and $S_n$ has the binomial distribution $B(n, p)$.

**¶31. Moment generating function.** The **moment generating function** of $X$ is defined to be

$$\phi_X(t) := \mathtt{E}[e^{tX}] = \sum_{i \geq 0} a_i e^{it}.$$

This is sometimes more convenient then the ordinary generating function. Differentiating $n$ times, we see $\phi_X^{(n)}(t) = \mathtt{E}[X^n e^{tX}]$ so $\phi^{(n)}(0)$ is the $n$th moment of $X$. For instance, if $X$ is $B(n, p)$ distributed then $\phi_X(t) = (pe^t + q)^n$.

_____EXERCISES

**Exercise 11.1:**
   (a) What is the generating function of the r.v. $X$ where $\{X = i\}$ is the event that a pair of independent dice roll yields a sum of $i$ ($i = 2, \ldots, 12$)?
   (b) What is the generating function of $c_0, c_1, \ldots$ where $c_i = 1$ for all $i$? Where $c_i = i$ for all $i$?                                                                    ◇


**Exercise 11.2:** Determine the generating functions of the following probability distributions:
   (a) binomial, (b) geometric, (c) poisson.                                            ◇


**Exercise 11.3:** Compute the mean and variance of the binomial distributed, exponential distributed and Poisson distributed r.v.'s using generating functions.                        ◇


_____END EXERCISES


## §12. Simple Randomized Algorithms


   Probabilistic thinking turns out to be uncannily effective for proving the existence of combinatorial objects. Such proofs can often be converted into randomized algorithms. There exist efficient randomized algorithms for problems that are not known to have efficient deterministic solutions. Even when both deterministic as well as randomized algorithms are available for a problem, the randomized algorithm is usually simpler. This fact may be enough to favor the randomized algorithm in practice. Sometimes, the route to deterministic solution is via a randomized one: after a randomized algorithm has been discovered, we may be able to remove the use of randomness. We will illustrate such "derandomization" techniques. For further reference, see Alon, Spencer and Erdös [1].

   We use a toy problem to illustrate probabilistic thinking in algorithm design. Suppose we want to color the edges of $K_n$, the complete graph on $n$ vertices with one of two colors, either red or blue. For instance, we can color all the edges red. But this is not a good choice if our goal is to be "as equitable as possible" evenly. Let us take this to mean that we want to have as few monochromatic triangles as possible. A triangle $(i, j, k)$ is monochromatic if all its edges have the same color. For instance, Figure 2(a,b) show two 2-colorings of $K_4$, having (a) one and (b) zero monochromatic triangles, respectively. We regard (b) are more equitable. In general, we cannot expect no monochromatic triangles, and we want to minimize this number.
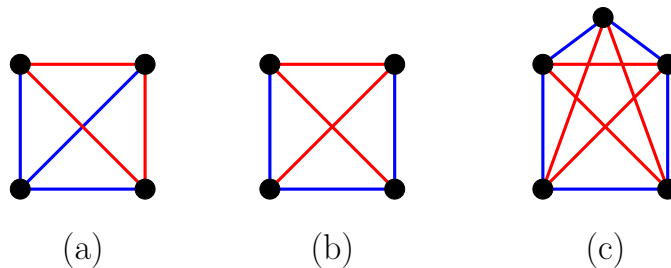


Figure 2: 2-colorings

   A *k-coloring* of the edges of a graph $G = (V, E)$ is an assignment $C : E \to \{1, \ldots, k\}$. There are $k^{|E|}$ such colorings. A *random k-coloring* is one that is equal to any of the $k^{|E|}$ colorings

with equal probability. Alternatively, a random $k$-coloring is one that assigns each edge to any of the $k$ colors with probability $1/k$. When $k = 2$, we assume the 2 colors are blue and red.

LEMMA 9. *Let $c$ and $n$ be positive integers. In the random 2-coloring of the edges of the complete graph $K_n$, the expected number of copies of $K_c$ that are monochromatic is*

$$\binom{n}{c} 2^{1-\binom{c}{2}}$$

*Proof.* Let $X$ count the number of monochromatic copies of $K_c$ in a random edge 2-coloring of $K_n$. If $V$ is the vertex set of $K_n$ then

$$X = \sum_U X_U \tag{42}$$

where $U \in \binom{V}{c}$ and $X_U$ is the indicator function of the event that the subgraph of $K_n$ restricted to $U$ is monochromatic. The probability that $X_U = 1$ is the probability that $U$ is monochromatic. Since there are $\binom{c}{2}$ edges in $U$ and we can color it monochromatic in one of two colors (all blue or all red), the probability of the event $X_U = 1$ is

$$\Pr\{X_U = 1\} = 2^{1-\binom{c}{2}}.$$

But $\mathrm{E}[X_U] = \Pr\{X_U = 1\}$ since $X_U$ is an indicator function. Since there are $\binom{n}{c}$ choices of $U$, we obtain $\mathrm{E}[X] = \sum_U \mathrm{E}[X_U] = \binom{n}{2} 2^{1-\binom{c}{2}}$, by linearity of expectation.      **Q.E.D.**

By one of the remarkable properties of expectation (§VIII.5) we conclude:

**Corollary 10.** *There exists an edge 2-coloring of $K_n$ such that the number of monochromatic copies of $K_c$ is at most*

$$\binom{n}{c} 2^{1-\binom{c}{2}}$$

**¶32. Simple Monte Carlo and Las Vegas Algorithms for good colorings.** Consider how to find a 2-coloring assured by this corollary. Furthermore, suppose $c = 3$ so that we desire a 2-coloring of $K_n$ such that the number of monochromatic triangles is at most $\binom{n}{3} 2^{1-\binom{3}{2}} = \frac{1}{4}\binom{n}{3}$. There is a trivial randomized algorithm if we are willing to settle for a slightly larger number of monochromatic triangles: let us say that a 2-coloring is **good** if it has at most $\frac{1}{3}\binom{n}{3}$ monochromatic triangles; otherwise it is **bad**.

---

RANDOMIZED COLORING ALGORITHM:
Input: $n$
Output: a good 2-coloring of $K_n$
     repeat forever:
        1.   Randomly color the edges of $K_n$ blue or red.
        2.   Count the number $X$ of monochromatic triangles.
        3.   If $X < \frac{1}{3}\binom{n}{3}$, return the random coloring.

---

If the program halts, the random coloring has the desired property. The probability that a random coloring is good is at least $1/4$. In proof: if the probability to be bad is more than $3/4$, then the expected number of monochromatic triangles in a random coloring is more than

---

$\frac{3}{4} \cdot \frac{1}{3} \binom{n}{3}$, contradicting our lemma. (This is just Markov's inequality.) Hence the probability of repeating the loop *at least once* is at most $3/4$. If $T$ is the time to do the loop once and $\widetilde{T}$ is the expected time of the algorithm, then

$$\widetilde{T} \le T + \frac{3}{4}\widetilde{T}$$

which implies $\widetilde{T} \le 4T$. But note that

$$T = T(n) = \mathcal{O}(n^3),$$

since there are $\mathcal{O}(n^3)$ triangles to check. Thus the expected running time is $\widetilde{T} = \mathcal{O}(n^3)$.     *a Las Vegas Algorithm!*

Note that our randomized algorithm has unbounded worst-case running time. Nevertheless, the probability that the algorithm halts is 1. Otherwise, if there is a positive probability $\epsilon > 0$ of not halting, then expected running time becomes unbounded ($\ge \epsilon \times \infty$), which is a contradiction. Alternatively, the probability of not halting is at most $(3/4)^\infty = 0$.

The above algorithm always give the correct answer and has *expected* polynomial running time. But if you are willing to accept a non-zero probability of error, you can get a *worst case* polynomial time. Here is such an algorithm: just return the first random coloring you get!    *a Monte Carlo Algorithm!* Clearly, it is now worst case $O(n^3)$ but there is a $3/4$ probability that the output is wrong (i.e., it has more than $\frac{1}{3}\binom{n}{3}$ monochromatic triangles. If you don't like this probability of error, you can reduce it to any $\epsilon > 0$ that you like. To get an error probability of at most $\epsilon = 0.001$, we just have to repeat the loop at most 25 times checking if the output is correct after each loop. The probability that you fail after 25 tries is less than 0.001. At that point, you can return the best coloring you have found so far.

--------------------------------------------------------------------------------Exercises

**Exercise 12.1:** Is it true that every 2-coloring of $K_6$ has some monochromatic triangle?    $\diamond$

**Exercise 12.2:** For small values of $n$, it seems easy to find 2-colorings with fewer than $\frac{1}{4}\binom{n}{3}$ monochromatic triangles. For instance, when $n = 5$, we can color any 5-cycle red and the non-cycle edges blue, then there are no monochromatic triangles (the theorem only gave a bound o 2 monochromatic triangles). Consider the following simple deterministic algorithm to 2-color $K_n$: pick any $T$ tour of $K_n$. A tour is an $n$-cycle that visits every vertex of $K_n$ exactly once and returns to the starting point. Color the $n$ edges in $T$ red, and the rest blue. Prove that this algorithm does not guarantee at most $\frac{1}{4}\binom{n}{3}$ monochrome triangles. For which values of $n$ does this algorithm give fewer than $\frac{1}{4}\binom{n}{3}$ monochrome triangles?    $\diamond$

**Exercise 12.3:** Fixed a $G$ graph with $n$ nodes. Show that a random graph of size $2\log n$ does not occur as an induced subgraph of $G$.    $\diamond$

**Exercise 12.4:**
(i) What is the role of "3/4" in the first randomized algorithm?
(ii) Give another proof that the probability of halting is 1, by lower bounding the probability of halting at the $i$th iteration.
(iii) Modify the algorithm into one that has a probability $\varepsilon > 0$ of not finding the desired coloring, and whose worst case running time is $\mathcal{O}_\varepsilon(n^3)$.
(v) Can you improve $T(n)$ to $o(n^3)$?    $\diamond$

**Exercise 12.5:**

(a) Construct a deterministic algorithm to 2-color $K_n$ so that there are at most $2\binom{n/2}{3}$ monochromatic triangles. HINT: use divide and conquer.

(b) Generalize this construction to giving a bound on the number of monochromatic $K_c$ for any constant $c \geq 3$. Compare this bound with the original probabilistic bound.    ◇

**Exercise 12.6:** Let $m, n, a, b$ be positive integers.

(a) Show that in a random 2-coloring of the edges of the complete bipartite graph $K_{m,n}$, the expected number of copies of $K_{a,b}$ that are monochromatic is

$$C(m, n, a, b) = \binom{m}{a}\binom{n}{b}2^{1-ab} + \binom{m}{b}\binom{n}{a}2^{1-ab}.$$

(b) Give a polynomial-time randomized algorithm to 2-color the edges of $K_{m,n}$ so that there are at most $C(m, n, a, b)/2$ copies of monochromatic $K_{2,2}$. What is the expected running time of your algorithm?    ◇

_____END EXERCISES

## §13.  Derandomization

The idea of derandomization is to give a deterministic algorithm modeled after a known randomized algorithm. We continue to use our randomized algorithm for finding good 2-colorings of $K_n$. The brute force way to do this is to systematically check each 2-colorings of $K_n$, and for each, check whether the coloring is good, i.e., the number of monochromatic triangles is $\leq \frac{1}{3}\binom{n}{3}$. Since there are $2^{\binom{n}{2}}$ such colorings, the complexity of this method is $\Omega(2^{\binom{n}{2}})$. Can we do better?

Let us step back a moment, to see the probabilistic setting of this brute force algorithm. We are really searching through the sample space $\Omega$ where $\Omega = \{0, 1\}^{|E_n|}$ where $E_n = \binom{V_n}{2}$ is the edge set of $K_n$ and $V_n = \{1, \ldots, n\}$. The exponential behaviour comes from the fact that $|\Omega|$ is exponential. But if we have a polynomial size sample space $\overline{\Omega}$, then the same brute force search becomes polynomial-time. Let us see what we need in general in our sample space.

- We want $(\overline{\Omega}, 2^{\overline{\Omega}}, \Pr)$ to be a discrete and uniform probability space.

- We must define in this probabilistic space an ensemble of $\binom{n}{2}$ indicator variables $X_{ij}$ $(1 \leq i < j \leq n)$. Each $X_{ij} = 1$ iff $i-j \in \binom{V_n}{2}$ is colored blue.

- 

- The ensemble must be 3-wise independent. This implies that for each triangle $\{i, j, k\} \in \binom{V_n}{3}$, the coloring of any edge in $\{i, j, k\}$ are independent of the other edges.

- Given $\omega \in \overline{\overline{\Omega}}$, we can determine the value of $X_{ij}(\omega)$ and hence count the number of monochromatic triangles in $\omega$.

_____

With this property, our algorithm goes as follows: for each $\omega \in \overline{\overline{\Omega}}$, we compute $X_{ij}(\omega)$ ($1 \le i \le j \le n$) in $O(n^2)$ time. For each triangle $(i, j, k)$, we determine if $X_{ij}(\omega) = X_{jk}(\omega) = X_{k,i}(\omega)$. If so, the triangle is monochromatic. Checking this for all triangles takes $O(n^3)$ time. We now know whether the coloring is good. If so, we stop and output the coloring. Otherwise, we try another $\omega$. This algorithm runs in time $O(n^3 |\overline{\Omega}|)$. Note that we can compute the optimal coloring (i.e., minimal number of monochromatic triangles) in $O(n^3 |\overline{\Omega}|)$ time.

**¶33. Constructing a 3-wise Independent Ensemble over a Small Sample Space.** We next address the problem of constructing $\overline{\Omega}$.

TO BE PROVI

**¶34. Arithmetic on Finite Fields.** Every finite field has size $p^k$ for some prime $p$ and $k \ge 1$. Up to isomorphism, these fields are uniquely determined by $p^k$ and so they are commonly denoted $GF(p^k)$. The case $k = 1$ is easy: $GF(p)$ is just $\mathbb{Z}_p = \{0, 1, \ldots, p-1\}$ with arithmetic operations modulo $p$. For $k \ge 2$, the elements of $GF(p^k)$ may be taken to be $GF(p)^k$, a $k$ vector in $GF(p)$. But arithmetic on $GF(p)^k$ is not the obvious one. The obvious attempt to define a field operations $\circ'$ in $GF(p)^k$ is by componentwise operations in $GF(p)$: $(x_1, \ldots, x_k) \circ' (y_1, \ldots, y_k) = (x_1 \circ y_1, \ldots, x_k \circ y_k)$. where $\circ$ is the corresponding field operation in $GF(p)$. Unfortunately, this does not result in a field. For instance, any element $(x_1, \ldots, x_k)$ which has a 0-component would become a zero-divisor. The correct way to define the field operations in $GF(p)^k$ is to take an irreducible polynomial $I(X)$ of degree $k$ and coefficients in $GF(p)$. Each element of $GF(p)^k$ is interpreted as a polynomial in $X$ of degree $k-1$. Now the field operations are taken to be corresponding field operation on polynomials, but modulo $I(X)$.

*GF stands for Galois Field*

Take the simplest case of $p = 2$. Suppose $k = 3$ and let $I(X) = X^3 + X + 1$. Each $(a, b, c) \in GF(2)^3$ is viewed as a polynomial $aX^2 + bX + c$. Addition is performed component-wise: $(a, b, c) + (a', b', c') = (a + a', b + b', c + c')$. But remember $a + a'$ is arithmetic modulo 2: in particular, $1 + 1 = 0$ and $-a = a$. But for multiplication, we may perform the ordinary polynomial multiplcation followed by reduction modulo $I(X)$. For instance, $X^3$ is reduced $X + 1$ (as $X^3 + X + 1 \equiv 0$ implies $X^3 \equiv -(X+1) \equiv X + 1$. Likewise, $X^4 \equiv X(X^3) \equiv X(X+1) = X^2 + X$. Thus the product $(a, b, c) \cdot (a', b', c')$ results in a polynomial

$$(aa')X^4 + (ab' + ba')X^3 + (ac' + bb' + ca')X^2 + (bc' + b'c)X + (cc')$$

which reduces to a polynomial of degree $\le 2$,

$$aa'(X^2 + X) + (ab' + a'b)(X + 1) + (ac' + bb' + ca')X^2 + (bc' + b'c)X + cc'$$

E.g., Let us compute the product $(1, 0, 1) \cdot (1, 1, 1)$ using the latter formula. We have

$$(1, 0, 1) \cdot (1, 1, 1) = (1, 1, 0) + (0, 1, 1) + 0 + (0, 1, 0) + (0, 0, 1) = (1, 1, 0).$$

Since $GF(p^k)$ is a field, we can do division. This can be reduced to computing multiplicative inverses. Given $P(X) \equiv 0$, suppose its inverse is $Q(X)$, i.e., $P(X)Q(X) \equiv 1$ modulo $I(X)$. This means

$$P(X)Q(X) + I(X)J(X) = 1$$

for some $J(X)$. We view this equation for polynomials in $GF(2)[X]$. Since $I(X)$ is irreducible and $P(X)$ is not a multiple of $I(X)$, the GCD of $I(X)$ and $P(X)$ is 1. Students may recall that if we compute the GCD of $P(X)$ and $I(X)$ using the "extended GCD" algorithm will produce the polynomials $Q(X)$ and $J(X)$. For instance, dividing $I(X) = X^3 + X + 1$ by

$P(X) = X + 1$ over $GF(2)[X]$, we obtain the quotient $Q(X) = X^2 + X$ with remainder of 1. Thus $P(X) \cdot Q(X) = 1 (\mathrm{mod}\, I(X))$. Hence the inverse of $P(X)$ is $Q(X) = X^2 + X$. In terms of bit vectors, this is the relation

$$(0, 1, 1) \cdot (1, 1, 0) = (0, 0, 1).$$

In the table below, we list the polynomials in $GF(2)[X]$ of degrees up to 3. For each polynomial, we either indicate that it is irreducible, or give its factorization. Thus $p_7 = X^2 + X + 1$ is the only irreducible polynomial of degree 2.

| Name | Polynomial | Irreducible? |
|------|-----------|--------------|
| $p_0$ | $0$ | ✓ |
| $p_1$ | $1$ | ✓ |
| $p_2$ | $X$ | ✓ |
| $p_3$ | $X + 1$ | ✓ |
| $p_4$ | $X^2$ | $p_0^2$ |
| $p_5$ | $X^2 + 1$ | $p_3^2$ |
| $p_6$ | $X^2 + X$ | $p_2 p_3$ |
| $p_7$ | $X^2 + X + 1$ | ✓ |
| $p_8$ | $X^3$ | $p_2 p_4$ |
| $p_9$ | $X^3 + 1$ | $p_3 p_7$ |
| $p_{10}$ | $X^3 + X$ | $p_2 p_5$ |
| $p_{11}$ | $X^3 + X^2$ | $p_2 p_6$ |
| $p_{12}$ | $X^3 + X + 1$ | ✓ |
| $p_{13}$ | $X^3 + X^2 + 1$ | ✓ |
| $p_{14}$ | $X^3 + X^2 + X$ | $p_2 p_7$ |
| $p_{16}$ | $X^3 + X^2 + X + 1$ | $p_3 p_5$ |

Table 1: Irreducible polynomials up to degree 3 over $GF(2)$

Modulo an irreducible polynomial $I(X)$ of degree $k$, we can construct a multiplication table for elements $GF(p^k)$, and list their inverses.

—————————————————————————————————————————————EXERCISES

**Exercise 13.1:** Give the multiplication and inverse tables for $GF(2^3)$ using $I(X) = p_7$.     ◇

**Exercise 13.2:** Implement the extended Euclidean algorithm for computing inverses in $GF(2^k)$ modulo an irreducible polynomial $I(X)$ of degree $k$.     ◇

**Exercise 13.3:** Although division can be reduced to inverses and multiplication, it would be more efficient to do this directly. Please adapt the extended GCD algorithm to produce such a direct division algorithm.     ◇

—————————————————————————————————————————————END EXERCISES

## §14. Tracking a Random Object

We began with an existence proof of a coloring of $K_n$ that does not have "too many" monochromatic copies of $K_c$. Then we derived a simple randomized algorithm to find such a coloring. Suppose we now want a deterministic algorithm instead. We use a method of Spencer and Raghavan to convert a randomized algorithm into a deterministic one. Alternatively, the method converts a probabilistic existence proof into a deterministic algorithm. Such a conversion has been termed "derandomization", and is based on conditional probabilities.

For our problem of 2-coloring $K_n$, the derandomization method is rather simple. We claim that the following deterministic algorithm will compute a 2-coloring with at most $\frac{1}{4}\binom{n}{3}$ monochromatic triangles:

---

Deterministic Coloring Algorithm:
1. Arbitrarily order the edges $e_1, e_2, \ldots, e_m$, $m = \binom{n}{2}$.
2. Consider each $e_i$ in turn:
     2.0. So $e_1, \ldots, e_{i-1}$ had been colored.
     2.1. Compute $W_i^{red}$, defined to be the expected number of monochromatic triangles
         if $e_i$ is next colored red and the remaining edges are randomly colored.
         Similarly, compute $W_i^{blue}$.
     2.2. Color $e_i$ red iff $W_i^{red} \leq W_i^{blue}$.

---

*randomized greediness!*

**¶35. Correctness.** We claim that the final coloring has at most $\frac{1}{4}\binom{n}{3}$ monochromatic triangles. In proof, let $W_i$ be the expected number of monochromatic triangles if the remaining edges $e_i, e_{i+1}, \ldots, e_m$ are randomly colored, conditioned a given coloring of $e_1, \ldots, e_{i-1}$. Initially we have

$$W_1 = \frac{1}{4}\binom{n}{3}.$$

Clearly,

$$W_i = \frac{W_i^{red} + W_i^{blue}}{2} \geq \min\{W_i^{red}, W_i^{blue}\} = W_i.$$

It follows that $W_1 \geq W_2 \geq \cdots \geq W_m$. But $W_m$ corresponds to a 2-coloring $C$ of $K_n$ and so $W_m$ must count the number of monochromatic triangles of $C$. Moreover,

$$W_m \leq W_0 \leq \frac{1}{4}\binom{n}{3},$$

as desired.

**¶36. Complexity.** To see that the above outline can be turned into an effective algorithm, we show that $W_i$ can be computed in polynomial time. This is straightforward: for any triple $U$ of vertices in $K_n$, let $X_{i,U}$ be the indicator function for the event that $U$ will be monochromatic if the remaining edges $e_{i+1}, \ldots, e_m$ are randomly colored. Then

$$W_i = \sum_U \mathrm{E}[X_{i,U}]$$

where the sum ranges over all $U$. But $\mathrm{E}[X_{i,U}]$ is just the probability that $U$ will become monochromatic:

$$\mathrm{E}[X_{i,U}] = \begin{cases} 2 \cdot 2^{-3} & \text{if no edge of } U \text{ has been colored,} \\ 2^{-3+i} & \text{if } i = 1, 2, 3 \text{ edges of } U \text{ has been colored with one color,} \\ 0 & \text{the edges of } U \text{ have been given both colors.} \end{cases}$$

---

Clearly, each $W_i^{red}$ and $W_i^{blue}$ can be computed in $\mathcal{O}(n^3)$ time. This leads to an $\mathcal{O}(n^5)$ time algorithm.

But we can improve this algorithm by a slight reorganization of the algorithm. Initially, compute $W_1$ in $O(n^3)$ time by summing over all triangles. Then for $i \geq 2$, we can compute $W_i^{red}$ and $W_i^{blue}$ in constant time by using the known value of $W_{i-1}$. Moreover, $W_i$ is updated to either $W_i^{red}$ or $W_i^{blue}$. Thus, the iterative loop takes $\mathcal{O}(n^2)$ time. The overall time is $O(n^3)$.

**¶37. Framework for deterministic tracking.** It is instructive to see the above algorithm in a general framework. Let $D$ be a set of objects and $\chi : D \to \mathbb{R}$ is a real function. For instance, $D$ is the set of 2-colorings of $K_n$ and $\chi$ counts the number of monochromatic triangles. In general, think of $\chi(d)$ as computing some statistic of $d \in D$. Call an object $d$ "good" if $\chi(d) \leq k$ (for some $k$); and "bad" otherwise. Our problem is to find a good object $d \in D$. First we introduce a sequence *tracking variables* $X_1, \ldots, X_m$ in some probability space $(\Omega, 2^\Omega, \mathrm{Pr})$. Each $X_i$ is an independent Bernoulli r.v. where $\mathrm{Pr}\{X_i = +1\} = \mathrm{Pr}\{X_i = -1\} = \frac{1}{2}$. We want these variables to be "complete" the sense that for any $\epsilon_1, \ldots, \epsilon_m \in \{\pm 1\}$, the event

$$\{X_1 = \epsilon_1, \ldots, X_m = \epsilon_m\}$$

is an elementary event. For instance, we can simply let $\Omega = \{\pm 1\}^m$ and

$$X_i(\epsilon_1, \epsilon_2, \ldots, \epsilon_m) = \epsilon_i$$

for $(\epsilon_1, \ldots, \epsilon_m) \in \Omega$. We also introduce a random object $g : \Omega \to D$ with the property

$$\mathrm{E}[\chi_g] \leq k \tag{43}$$

where $\chi_g$ is the random variable $\chi_g(\omega) = \chi(g(\omega))$, $\omega \in \Omega$. This means that there is *some* sample point $\omega$ such that $g(\omega)$ is good. Write $W_i$ as shorthand for $W(\epsilon_1, \ldots, \epsilon_{i-1})$ where $W(\epsilon_1, \ldots, \epsilon_{i-1})$ is the conditional expectation

$$W(\epsilon_1, \ldots, \epsilon_{i-1}) := \mathrm{E}[\chi_g | X_1 = \epsilon_1, \ldots, X_i = \epsilon_i].$$

Hence, our assumption (43) above amounts to $W_0 \leq k$. Inductively, suppose we have determined $\epsilon_1, \ldots, \epsilon_{i-1}$ so that

$$W_{i-1} = W(\epsilon_1, \ldots, \epsilon_{i-1}) \leq k.$$

Then observe that

$$
\begin{aligned}
W_{i-1} &= \frac{W(\epsilon_1, \ldots, \epsilon_{i-1}, +1) + W(\epsilon_1, \ldots, \epsilon_{i-1}, -1)}{2} \\
&\geq \min\{W(\epsilon_1, \ldots, \epsilon_{i-1}, +1), W(\epsilon_1, \ldots, \epsilon_{i-1}, -1)\}.
\end{aligned}
$$

Hence, if we can efficiently compute the value $W(\epsilon_1', \ldots, \epsilon_i')$ for any choice of $\epsilon_j'$'s, we may choose $\epsilon_i$ so that $W_i \leq W_{i-1}$, thus extending our inductive hypothesis. The hypothesis $W_i \leq k$ implies there is a sample point $\omega \in \{X_1 = \epsilon_1, \ldots, X_i = \epsilon_i\}$ such that $g(\omega)$ is good. In particular, the inequality $W_m \leq k$ implies that $g(\omega)$ is good, where $\omega = (\epsilon_1, \ldots, \epsilon_m)$. Thus, after $m$ steps, we have successfully "tracked" down a good object $g(\omega)$. Note that this is reminiscent of the greedy approach.

The use of the conditional expectations $W_i$ is clearly central to this method. A special case of conditional expectation is when $W_i$ are conditional probabilities. If we cannot efficiently compute the $W_i$'s, some estimates must be used. This will be our next illustration.

_____Exercises

**Exercise 14.1:** Generalize the above derandomized algorithm 2-coloring $K_n$ while avoiding too many monochromatic $K_c$, for any $c \geq 4$. What is the complexity of the algorithm?

$\Diamond$

**Exercise 14.2:**

(i) If the edges of $K_4$ (the complete graph on 4 vertices) are 2-colored so that two edges $e, e'$ have one color and the other 4 edges have a different color, and moreover $e, e'$ have no vertices in common, then we call this coloring of $K_4$ a "kite". What is the expected
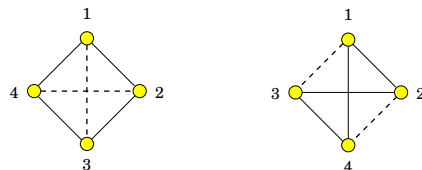


Figure 3: A kite drawn in two ways: edges (1,3) and (2,4) are red, the rest blue.

number of kites in a random 2-coloring of the edges of $K_n$?
(ii) Devise a deterministic tracking algorithm to compute a 2-coloring which achieves *at least* this expected number. but analyze its complexity.
(iii) Forget deterministic tracking, but give a simple $\mathcal{O}(n^2)$ algorithm to do the same. $\Diamond$

_____END EXERCISES

## §15. Maximum Satisfiability

In the classic Satisfiability Problem (SAT), we are given a Boolean formula $F$ over some set $\{x_1, \ldots, x_n\}$ of Boolean variables, and we have to decide if $F$ is satisfiable by some assignment of truth values to the $x_i$'s. We may assume $F$ is given as the conjunction of of clauses $F = \bigwedge_{i=1}^m C_i$ ($m \geq 1$) and each clause $C_i$ is a disjunction of one or more Boolean literals (a literal is either a variable $x_i$ or its negation $\overline{x_i}$). As usual, we can view $C_i$ as a set of Boolean literals and $F$ as a set of clauses. For instance, let

$$F_0 = \{C_1, C_2, C_3, C_4\}, \quad C_1 = \{\overline{x}\}, C_2 = \{x, \overline{y}\}, C_3 = \{y, \overline{z}\}, C_4 = \{x, y, z\} \tag{44}$$

where the Boolean variables are $x, y, z$. It is easy to see that $F_0$ is not satisfiable here.

The **Maximum Satisfiability Problem** (MAXSAT) is variant of SAT in which we just want to maximize the number of satisfied clauses in $F$. It is well-known that SAT is $NP$-hard, and clearly MAXSAT remains $NP$-hard. But MAXSAT is now amenable to an approximation interpretation: to find an assignment that satisfies at least a constant fraction of the optimal number of satisfiable clauses. For the unsatisfiable formula $F_0$ in (44), we see that we can satisfy 3 out of the four clauses by setting $x, y, z$ to 1.

If $A$ is an algorithm for MAXSAT, let $A(F)$ denote the number of clauses satisfied by this algorithm on input $F$. Also, let $A^*(F)$ denote the maximum satisfiable number of clauses for input $F$ (think of $A^*$ as the optimal algorithm). Fix a constant $0 \leq \alpha \leq 1$. We call $A$ an $\alpha$-**approximation algorithm** for MAXSAT if $A(F)/A^*(F) \geq \alpha$ for all input $F$. This definition

extends to the case where $A$ is a randomized algorithm. In this case, we replace $A(F)$ by its expected value $\mathbf{E}[A(F)]$.

This approximate MAXSAT problem is also amenable to our tracking framework. But in this section, we want to illustrate another technique: **random rounding**. The idea is to first give an optimal "fractional solution" in which each Boolean variable are initially assigned a fraction between 0 and 1. We then randomly round these fractional values to 0 or 1; we would expect the rounded result to achieve an provably good objective value. Note that is a randomized algorithm.

**¶38. Initial Algorithm $A_0$.** Before we proceed deeper, consider the following trivial algorithm $A_0$: on input $F = \{C_1, \ldots, C_m\}$, we set each Boolean variable to 0 or 1 with equal probability. Let $X_i$ be the indicator r.v. for the event that $C_i$ is satisfied. Then $\mathbf{E}[X_i] = \Pr\{X_i = 1\} = 1 - 2^{-k}$ where $k$ is the number of literals in $C_i$. Since $k \geq 1$, we get $\mathbf{E}[X_i] \geq 1/2$. The expected number of clauses satisfied is therefore equal to

$$\mathbf{E}[A_0(F)] = \sum_{i=1}^{m} \mathbf{E}[X_i] = m/2.$$

This proves that $A_0$ is a 1/2-approximation algorithm for MAXSAT (since $A^*(F) \leq m$). Of course, the bound $\alpha = 1/2$ can be improved if we know that each clause has more than 1 literal. For instance, if $|C_i| \geq 2$, then $A_0$ is already a 3/4-approximation algorithm for MAXSAT.

We now return to our goal of $\alpha = 3/4$, without assumptions such as $|C_i| \geq 2$. The idea is to turn the satisfiability problem into a linear programming problem. This transformation is well-known: We introduce the real variables $c_1, \ldots, c_m$ ($c_i$ corresponding to clause $C_i$), and $v_1, \ldots, v_n$ ($v_j$ corresponding to variable $x_j$). Intuitively, $v_j$ is the truth value assigned to $x_j$, and $c_i = 1$ if $C_i$ is satisfied. However, as a linear programming problem, we can only write write a set of linear constraints on these variables:

$$0 \leq v_j \leq 1, \quad 0 \leq c_i \leq 1.$$

Although we intend the $v_j$'s and $c_i$'s to be either 0 or 1, these inequalities only constrain them to lie between 0 and 1. To connect these two set of variables, we look at each clause $C_i$ and write the inequality

$$c_i \leq \left( \sum_{j:x_j \in C_i} v_j \right) + \left( \sum_{j:\overline{x_j} \in C_i} (1 - v_j) \right). \tag{45}$$

The objective of the linear programming problem is to maximize the sum

$$c_1 + \cdots + c_m \tag{46}$$

subject to the preceding inequalities. Let us verify that the inequalities (45) ($i = 1, \ldots, m$) would capture the MAXSAT problem correctly if the variables are truly $0-1$. Suppose $v_1, \ldots, v_n$ are $0-1$ values representing an assignment. If $C_i$ is satisfied by this assignment, then at least one of the two sums on the righthand side of (45) is at least 1, and so (45) is clearly satisfied. Moreover, we should set $c_i = 1$ in order to maximize the objective function (46). Conversely, if $C_i$ is not satisfied, then the righthand side of (45) is zero, and we must have $c_i = 0$. This proves that our linear program exactly captures the MAXSAT problem with the proviso *the variables $v_j$ assume integer values*.

Let $\widetilde{A}(F)$ denote the maximum value (46) of our linear program. Since our linear programming solution is maximized over possibly non-integer values of $v_j$'s (thus failing our proviso), we conclude that

$$\widetilde{A}(F) \geq A^*(F). \tag{47}$$

¶39. The Random Rounding Algorithm $A_1$. We now describe a **random rounding algorithm**, $A_1$: on input $F$, we set up the linear program corresponding to $F$, and solve it. Let this linear program have solution

$$(v_1', v_2', \ldots, v_n'; c_1', \ldots, c_m').$$

Thus each $v_j'$ and $c_i'$ is a value between 0 and 1. Clearly (45) implies:

$$c_i' = \min \left\{ 1, \left( \sum_{j:x_j \in C_i} v_j \right) + \left( \sum_{j:\overline{x_j} \in C_i} (1 - v_j) \right) \right\}. \tag{48}$$

Since the $v_j'$s lie between 0 and 1, we can interpret them as probabilities. Our algorithm will assign $x_j$ to the value 1 with probability $v_j'$, and to the value 0 otherwise. This completes our description of $A_1$.

Our algorithm $A_1$ calls some linear program solver as a subroutine. Polynomial time algorithms for such solvers are known[6] and are based on the famous interior-point method.

Let us now consider the $0-1$ solution $(x_1', \ldots, x_n')$ computed by $A_1$. What is the probability that it fails to satisfy a clause $C_i$? This probability is seen to be

$$p_i := \left( \prod_{j:x_j \in C_i} (1 - v_j') \right) \cdot \left( \prod_{j:\overline{x_j} \in C_i} v_j' \right).$$

Therefore the probability that $C_i$ is satisfied is $1 - p_i$. For simplicity in the notations, assume that $C_i$ has no negated variables and $C_i$ has $k$ variables. Then

$$p_i = \prod_{j:x_j \in C_i} (1 - v_j')$$

subject to the constraint

$$\sum_{j:x_j \in C_i} v_j' \geq z_i'.$$

We see that $p_i$ is maximized when $v_j' = z_i'/|C_i| = z_i'/k$. This proves that

$$p_i \leq (1 - z_i'/k)^k \leq (1 - 1/k)^k.$$

## §16. Discrepancy Random Variables

A typical "discrepancy problem" is this: *given real numbers $a_1, \ldots, a_n$, choose signs $\epsilon_1, \ldots, \epsilon_n \in \{\pm 1\}$ so as to minimize the absolute value of the sum*

$$S = \sum_{i=1}^{n} \epsilon_i a_i.$$

The minimum value of $|S|$ is the *discrepancy* of $(a_1, \ldots, a_n)$. Call[7] a random variable $X$ a *discrepancy r.v.* if the range of $X$ is $\pm 1$; it is *random* if, in addition,

$$\Pr\{X = +1\} = \Pr\{X = -1\} = \frac{1}{2}.$$

---

[6] Interestingly, if the number of variables is bounded by some constant, then linear programming has rather simple randomized algorithms. But in this application, we cannot bound the number of variables by a constant.
[7] Clearly, this is just another name for a Bernoulli r.v..

The *hyperbolic cosine function* $\cosh(x) = (e^x + e^{-x})/2$ arises naturally in discrepancy random variables. If $X_i$ are random discrepancy r.v.'s, then

$$
\begin{aligned}
\mathbb{E}[e^{a_i X_i}] &= \frac{e^{a_i} + e^{-a_i}}{2} \\
&= \cosh(a_i) \\
\mathbb{E}[e^{a_1 X_1 + a_2 X_2}] &= \frac{e^{a_1 + a_2} + e^{-a_1 - a_2} + e^{a_1 - a_2} + e^{-a_1 + a_2}}{4} \\
&= \frac{\cosh(a_1 + a_2) + \cosh(a_1 - a_2)}{2}.
\end{aligned}
$$

Using the fact that

$$
2\cosh(a_1)\cosh(a_2) = \cosh(a_1 + a_2) + \cosh(a_1 - a_2),
$$

we conclude that $\mathbb{E}[e^{a_1 X_1 + a_2 X_2}] = \cosh(a_1)\cosh(a_2)$. In general, with $S = \sum_{i=1}^{n} a_i X_i$, we get

$$
\mathbb{E}[e^S] = \prod_{i=1}^{n} \cosh(a_i). \tag{49}
$$

A useful inequality in this connection is

$$
\cosh(x) \leq e^{x^2/2}, \qquad x \in \mathbb{R}, \tag{50}
$$

with equality iff $x = 0$. This can be easily deduced from the standard power series for $e^x$.

_____Exercises

**Exercise 16.1:**
    (i) Verify equation (50).
    (ii) Show the bound $\Pr\{X_1 + \cdots + X_n > a\} < e^{-a^2/2n}$, where $X_i$ are random discrepancy r.v.'s and $a > 0$.       $\Diamond$

_____End Exercises

## §17. A Matrix Discrepancy Problem

    Raghavan considered a discrepancy problem in which we need to estimate the conditional probabilities. Let $A = (a_{ij})$ be an $n \times n$ input matrix with $|a_{ij}| \leq 1$. Let $\Omega = \{\pm 1\}^n$. Our goal want to find $\bar{\epsilon} = (\epsilon_1, \ldots, \epsilon_n) \in \Omega$, such that for each $i = 1, \ldots, n$,

$$
\left| \sum_{j=1}^{n} \epsilon_j a_{ij} \right| \leq \alpha n
$$

where

$$
\alpha := \sqrt{\frac{2\ln(2n)}{n}}. \tag{51}
$$

This choice of $\alpha$ will fall out from the method, so it is best to treat it as a yet-to-be-chosen constant ($n$ is fixed during this derivation). Using the method of deterministic tracking, we

introduce random discrepancy r.v.'s $X_1, \ldots, X_n$ such that $X_i(\bar{\epsilon}) = \epsilon_i$ for all $i$. Also introduce the r.v.'s

$$S_i = \sum_{j=1}^{n} X_j a_{ij}, \qquad i = 1, \ldots, n.$$

Suppose the values $\epsilon_1, \ldots, \epsilon_\ell$ have been choosen. Consider the event

$$C^\ell := \{X_1 = \epsilon_1, \ldots, X_\ell = \epsilon_\ell\}$$

and the conditional "bad" event

$$B_i^\ell := \{\,|S_i| > \alpha n \,|\, C^\ell\,\}.$$

To carry out the deterministic tracking method above, we would like to compute the probability $\Pr(B_i^\ell)$. Unfortunately we do not know how to do this efficiently. We therefore replace $\Pr(B_i^\ell)$ by an easy to compute upper estimate, as follows:

$$
\begin{aligned}
\Pr(B_i^\ell) &= \Pr\{|S_i| > \alpha n \,\big|\, C^\ell\} \\
&= \Pr\{e^{\alpha S_i} > e^{\alpha^2 n} \,\big|\, C^\ell\} + \Pr\{e^{-\alpha S_i} > e^{\alpha^2 n} \,\big|\, C^\ell\} \\
&\leq e^{-\alpha^2 n} \mathrm{E}[e^{\alpha S_i} + e^{-\alpha S_i} \,\big|\, C^\ell] \quad \text{(Markov inequality)} \\
&= e^{-\alpha^2 n} W_i^\ell,
\end{aligned}
$$

where the last equation defines $W_i^\ell$. Thus we use $W_i^\ell$ as surrogate for $\Pr(B_i^\ell)$. We do it because we can easily compute $W_i^\ell$ as follows:

$$
\begin{aligned}
W_i^\ell &= \mathrm{E}[e^{\alpha S_i} + e^{-\alpha S_i} | C^\ell] \\
&= \mathrm{E}[\exp\left(\alpha \sum_{j=1}^{\ell} \epsilon_j a_{ij}\right) \exp\left(\alpha \sum_{j=\ell+1}^{n} X_j a_{ij}\right)] + \mathrm{E}[\exp\left(-\alpha \sum_{j=1}^{\ell} \epsilon_j a_{ij}\right) \exp\left(-\alpha \sum_{j=\ell+1}^{n} X_j a_{ij}\right)] \\
&= \exp\left(\alpha \sum_{j=1}^{\ell} \epsilon_j a_{ij}\right) \prod_{j=\ell+1}^{n} \cosh(\alpha a_{ij}) + \exp\left(-\alpha \sum_{j=1}^{\ell} \epsilon_j a_{ij}\right) \prod_{j=\ell+1}^{n} \cosh(\alpha a_{ij}) \\
&= 2\cosh\left(\alpha \sum_{j=1}^{\ell} \epsilon_j a_{ij}\right) \prod_{j=\ell+1}^{n} \cosh(\alpha a_{ij}).
\end{aligned}
$$

In particular, for $\ell = 0$,

$$
\begin{aligned}
\sum_{i=1}^{n} W_i^0 &= 2 \sum_{i=1}^{n} \prod_{j=1}^{n} \cosh(\alpha a_{ij}) \\
&\leq 2 \sum_{i=1}^{n} \prod_{j=1}^{n} e^{a_{ij}^2 \alpha^2 / 2} \\
&< 2n e^{n\alpha^2/2},
\end{aligned}
$$

where the last inequality is strict since we will assume no row of $A$ is all zero. So the probability that a random choice of $\bar{\epsilon}$ is bad is at most

$$
\begin{aligned}
\sum_{i=1}^{n} \Pr(B_i^0) &\leq e^{-n\alpha^2} \sum_{i=1}^{n} W_i^0 \\
&< 2n e^{-n\alpha^2/2}.
\end{aligned}
$$

We choose $\alpha$ so that the last expression is equal to 1; this is precisely the $\alpha$ in (51). This proves that there is a sample point $\bar{\epsilon}$ where none of the $n$ bad events $B_i^0$ occur. The problem now is to

---

track down this sample point. In the usual fashion, we show that if $\epsilon_1, \ldots, \epsilon_\ell$ have been chosen then we can choose $\epsilon_{\ell+1}$ such that

$$\sum_{i=1}^{n} W_i^\ell \geq \sum_{i=1}^{n} W_i^{\ell+1}.$$

This can be done as follows: let $C_+^\ell$ denote the event $C^\ell \cap \{X_{\ell+1} = +1\}$ and similarly let $C_-^\ell := C^\ell \cap \{X_{\ell+1} = -1\}$. Then

$$
\begin{aligned}
\sum_{i=1}^{n} W_i^\ell &= \sum_{i=1}^{n} \mathbf{E}[e^{\alpha S_i} + e^{-\alpha S_i} | C^\ell] \\
&= \sum_{i=1}^{n} \frac{\mathbf{E}[e^{\alpha S_i} + e^{-\alpha S_i} | C_\ell^+] + \mathbf{E}[e^{\alpha S_i} + e^{-\alpha S_i} | C_\ell^-]}{2} \\
&\geq \min\{\sum_{i=1}^{n} \mathbf{E}[e^{\alpha S_i} + e^{-\alpha S_i} | C_\ell^+], \sum_{i=1}^{n} \mathbf{E}[e^{\alpha S_i} + e^{-\alpha S_i} | C_\ell^+]\} \\
&= \sum_{i=1}^{n} \mathbf{E}[e^{\alpha S_i} + e^{-\alpha S_i} | C^{\ell+1}]
\end{aligned}
$$

provided we choose $\epsilon_{\ell+1}$ to make the last equation hold. But the final expression (52) is $\sum_{i=1}^{n} W_i^{\ell+1}$. This proves $\sum_{i=1}^{n} W_i^\ell \geq \sum_{i=1}^{n} W_i^{\ell+1}$. After $n$ steps, we have

$$\sum_{i=1}^{n} \Pr(B_i^n) \leq e^{-\alpha^2 n} \sum_{i=1}^{n} W_i^n < 1. \tag{52}$$

But $B_i^n$ is the probability that $|S_i| > \alpha n$, conditioned on the event $C^n$. As $C^n = \{(\epsilon_1, \ldots, \epsilon_n)\}$ is an elementary event, the probability of any event conditioned on $C^n$ is either 0 or 1. Thus equation (52) implies that $\Pr(B_i^n) = 0$ for all $i$. Hence $C^n$ is a solution to the discrepancy problem.

We remark that computing $W_i$ is considered easy because the exponential function $e^x$ can be computed relatively efficiently to any desired degree of accuracy (see Exercise).

_____EXERCISES

**Exercise 17.1:** What is the bit-complexity of Raghavan's algorithm? Assume that $e^x$ (for $x$ in any fixed interval $[a, b]$) can be computed to $n$-bits of relative precision in $\mathcal{O}(M(n) \log n)$ time where $M(n)$ is the complexity of binary number multiplication. The inputs numbers $a_{ij}$ are in floating point notation, *i.e.*, $a_{ij}$ is represented as a pair $(e_{ij}, f_{ij})$ of binary integers so that

$$a_{ij} = 2^{e_{ij}} f_{ij}$$

and $e_{ij}, f_{ij}$ are at most $m$-bit numbers.                                                    ◇

_____END EXERCISES

## §18. Random Search Trees

Recall the concept of binary search trees from Lecture III. This lecture focuses on a class of random binary search trees called **random treaps**. Basically, treaps are search trees whose shape is determined by the assignment of **priorities** to each key. If the priorities are chosen randomly, the result is a random treap with expected height of $\Theta(\log n)$. Why is expected $\Theta(\log n)$ height interesting when worst case $\Theta(\log n)$ is achievable? One reason is that randomized algorithms are usually simpler to implement. Roughly speaking, they do not have to "work too hard" to achieve balance, but can rely on randomness as an ally.

The analysis of random binary search trees has a long history. The analysis of binary search trees under random insertions only was known for a long time. It had been an open problem to analyze the expected behavior of binary search trees under insertions *and deletions*. A moment reflection will indicate the difficulty of formulating such a model. As an indication of the state of affairs, a paper *A trivial algorithm whose analysis isn't* by Jonassen and Knuth [6] analyzed the random insertion and deletion of a binary tree containing no more than 3 items at any time. The currently accepted probabilistic model for random search trees is due to Aragon and Seidel [2] who introduced the treap data-structure.[8] Prior to the treap solution, Pugh [14] introduced a simple but important data structure called **skip list** whose analysis is similar to treaps. The real breakthrough of the treap solution lies in its introduction of a new model for randomized search trees, thus changing[9] the ground rules for analyzing random trees.

**Notations.** We write $[i..j]$ for the set $\{i, i+1, \ldots, j-1, j\}$ where $i \leq j$ are integers. Often, the set of keys is $[1..n]$. If $u$ is an item, we write $u.\texttt{key}$ and $u.\texttt{priority}$ for the key and priorities associated with $u$. Since we usually confuse an item with the node it is stored in, we may also write $u.\texttt{parent}$ or $u.\texttt{leftChild}$. By convention, that $u.\texttt{parent} = u$ if $u$ is the root and $u.\texttt{leftChild} = u$ if $u$ has no left child.

## §19. Skip Lists

It is instructive to first look at the skip list data structure. Pugh [14] gave experimental evidence that its performance is comparable to non-recursive AVL trees algorithms, and superior to splay trees (the experiments use $2^{16}$ integer keys).

Suppose we want to maintain an ordered list $L_0$ of keys $x_1 < x_2 < \cdots < x_n$. We shall construct a **hierarchy** of sublists

$$L_0 \supseteq L_1 \supseteq \cdots \supseteq L_{m-1} \supseteq L_m$$

where $L_m$ is only empty list in this hierarchy. Each $L_{i+1}$ is a *random sample* of $L_i$ in the following sense: each key of $L_i$ is put into $L_{i+1}$ with probability $1/2$. The hierarchy stops the first time the list becomes empty. In storing these lists, we shall add an artificial key $-\infty$ at the head of each list. Let $L'_i$ be the list $L_i$ augmented with $-\infty$. It is sufficient to store the lists $L'_i$ as a linked list with $-\infty$ as the head (a singly-linked list is sufficient). In addition, corresponding keys in two consecutive levels shares a two-way link (called the up- and down-links, respectively. This is illustrated in figure 4.

Clearly, the expected length of $L_i$ is $\mathrm{E}[|L_i|] = n2^{-i}$. Let us compute the expected height $\mathrm{E}[m]$. The probability of any key of $L_0$ appearing in level $i \geq 0$ is $1/2^i$. Now $m \geq i+1$ iff some

---

[8] The name "treap" comes from the fact that these are binary search trees with the heap property. It was originally coined by McCreight for a different data structure.

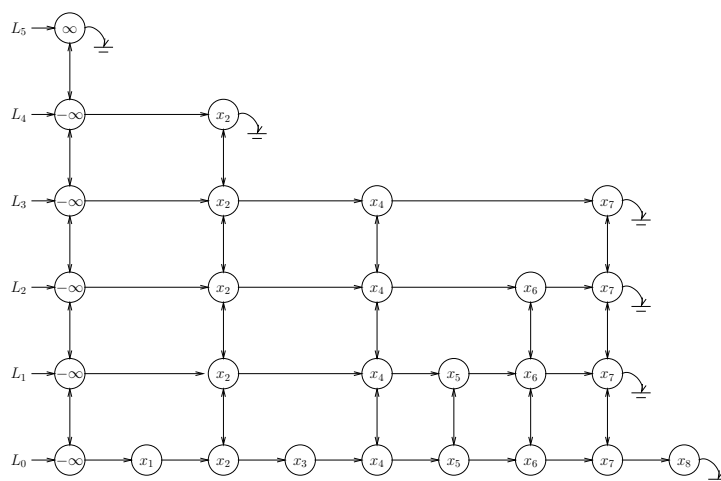[9] Another historical example in the same category is Alexander the Great's solution to the Gordan Knot problem.

Figure 4: A skip list on 8 keys.

key $x_j$ appears level $i$. Since there are $n$ choices for $j$, we have

$$\Pr\{m \geq i+1\} \leq \frac{n}{2^i}.$$

Hence

$$
\begin{aligned}
\mathrm{E}[m] &= \sum_{i \geq 1} i \Pr\{m = i\} \\
&= \sum_{i \geq 1} \Pr\{m \geq i\} && \text{(standard trick)} \\
&= \sum_{i \leq \lceil \lg n \rceil} \Pr\{m \geq i\} + \sum_{j > \lceil \lg n \rceil} \Pr\{m \geq j\} && \text{(another one!)} \\
&\leq \lceil \lg n \rceil + \sum_{j > \lceil \lg n \rceil} \frac{n}{2^{j-1}} \\
&\leq \lceil \lg n \rceil + \sum_{j > \lceil \lg n \rceil} \frac{1}{2^{j-1-\lg n}} \\
&\leq \lceil \lg n \rceil + 2.
\end{aligned}
$$

Have we really achieved anything special? You rightly think: why not just **deterministically** omit every other item in $L_i$ to form $L_{i+1}$? then $m \leq \lceil \lg n \rceil$ with less fuss! The reason why the probabilistic solution is superior is because the analysis will hold up even in the presence of insertion or deletion of arbitrary items. The deterministic solution may look very bad for particular sequence of deletions, for instance (how?). The critical idea is that the probabilistic behavior of each item $x$ is *independent* of the other items in the list: it is determined by its own sequence of coin tosses to decide whether to promote $x$ to the next level.

**Analysis of Lookup.**    Our algorithm to lookup a key $k$ in a skip list returns the largest key $k_0$ in $L_0$ such that $k_0 \leq k$. If $k$ is less than the smallest key $x_1$ in $L_0$, we return the special value $k_0 = -\infty$.

You may also be interested in finding the smallest key $k_1$ in $L_0$ such that $k_1 \geq k$. But $k_1$ is either $k_0$ or the successor of $k_0$ in $L_0$. If $k_0$ has no successor, then we imagine "$+\infty$" as its successor.

In general, define $k_i$ ($i = 0, \ldots, m$) to be the largest key in level $L'_i$ that is less than or equal to $k$. So $k_m = -\infty$. In general, given $k_{i+1}$ in level $i + 1$, we can find $k_i$ by following a down-link and then "scanning forward" in search of $k_i$. The search ends after the scan in list $L_0$: we end up with the element $k_0$. Let us call the sequence of nodes visited from $k_m$ to $k_0$

the **scan-path**. Let $s$ be the random variable denoting the length of the scan-path. If $s_i$ is the number of nodes of the scan-path that that lies in $L'_i$, then we have

$$s = \sum_{i=0}^{m} s_i.$$

The expected cost of looking up key $k$ is $\Theta(\mathbb{E}[s])$.

LEMMA 11. $\mathbb{E}[s] < 2\lceil \lg n \rceil + 4$.

Before presenting the proof, it is instructive to give a wrong argument: since $\mathbb{E}[m] \leq \ln n + 3$, and $\mathbb{E}[s_i] < 2$ (why?), we have $\mathbb{E}[s] \leq 2(\ln n + 2)$. This would be correct if each $s_i$ is i.i.d. and also independent of $m$ (see §VII.5). Unfortunately, $s_i$ is not independent of $m$ (e.g., $s_i = 0$ iff $m < i$).

*Proof of lemma 11.* Following Pugh, we introduce the random variable

$$s(\ell) = \sum_{i=0}^{\ell-1} s_i$$

where $\ell \geq 0$ is any constant (the threshhold level). By definition, $s(0) = 0$. We note that

$$s \leq s(\ell) + |L_\ell| + (m - \ell + 1).$$

To see this, note that the edges in the scan-path at levels $\ell$ and above involve at most $|L_\ell|$ horizontal edges, and at most $m - \ell$ vertical (down-link) edges.

Choose $\ell = \lceil \lg n \rceil$. Then $\mathbb{E}[|L_\ell|] = n/2^\ell \leq 1$. Also $\mathbb{E}[m - \ell + 1] \leq 3$ since $\mathbb{E}[m] \leq \lceil \lg n \rceil + 2$. The next lemma will imply $\mathbb{E}[s(\ell)] \leq 2\lceil \lg n \rceil$. Combining these bounds yields lemma 11.

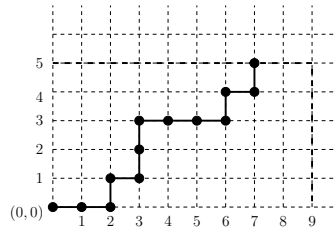**A Random Walk.** To bound $s(\ell)$, we introduce a random walk on the lattice $\mathbb{N} \times \mathbb{N}$.



Figure 5: Random walk with $X_{9,5} = 12$.

Let $u_t$ denote the position at time $t \in \mathbb{N}$. We begin with $u_0 = (0,0)$, and we terminate at time $t$ if $u_t = (x,y)$ with either $x = n$ or $y = \ell$. At each discrete time step, we can move from the current position $(x,y)$ to $(x,y+1)$ with probability $p$ $(0 < p < 1)$ or to $(x+1,y)$ with probability $q = 1 - p$. So $p$ denotes the probability of "promotion" to the next level. Let $X_{n,\ell}$ denoting the number of steps in this random walk at termination. For $p = 1/2$, we obtain the following connection with $s(\ell)$:

$$\mathbb{E}[X_{n,\ell}] \geq \mathbb{E}[s(\ell)]. \tag{53}$$

To see this, imagine walking along the scan-path of a lookup, but in the reverse direction. We start from key $k_0$. At each step, we get to move up to the next level with probability $p$ or to

scan backwards one step with probability $q = 1 - p$. There is a difference, however, between a horizontal step in our random walk and a backwards scan of the scan-path. In the random walk, each step is unit length, but the backwards scan may take a step of larger lengths with some non-zero probability. But this only means that $X_{n,\ell} \succeq s(\ell)$ (denoting stochastic domination, §VII.3). Hence (53) is an inequality instead of an equality. For example, in the skip list of figure 4, if $k_0 = x_8$ then $s(\ell) = s(4) = 7$ while $X_{n,\ell} = X_{8,4} = 10$.

LEMMA 12. *For all $\ell \geq 0$, $\mathrm{E}[X_{n,\ell}] \leq \ell/p$.*

*Proof.* Note that this bound removes the dependence on $n$. Clearly $\mathrm{E}[X_{n,0}] = 0$, so assume $\ell \geq 1$. Consider what happens after the first step: with probability $p$, the random walk moves to the next level, and the number of steps in rest of the random walk is distributed as $X_{n,\ell-1}$; similarly, with probability $1 - p$, the random walk remain in the same level and the number of steps in rest of the random walk is distributed as $X_{n-1,\ell}$. This gives the recurrence

$$\mathrm{E}[X_{n,\ell}] = 1 + p \cdot \mathrm{E}[X_{n,\ell-1}] + q \cdot \mathrm{E}[X_{n-1,\ell}].$$

Since $\mathrm{E}[X_{n-1,\ell}] \leq \mathrm{E}[X_{n,\ell}]$, we have

$$\mathrm{E}[X_{n,\ell}] \leq 1 + p \cdot \mathrm{E}[X_{n,\ell-1}] + q \cdot \mathrm{E}[X_{n,\ell}].$$

This simplifies to

$$\mathrm{E}[X_{n,\ell}] \leq \frac{1}{p} + \mathrm{E}[X_{n,\ell-1}].$$

This implies $\mathrm{E}[X_{n,\ell}] \leq \frac{\ell}{p}$ since $\mathrm{E}[X_{n,0}] = 0$.                    **Q.E.D.**

If $\ell = \lceil \lg n \rceil$ and $p = 1/2$ then

$$\mathrm{E}[s(\ell)] \leq \mathrm{E}[X_{n,\ell}] \leq 2 \lceil \lg n \rceil,$$

as noted in the proof of lemma 11.

We leave as an exercise to the reader to specify, and to analyze, algorithms for insertion and deletion in skip lists.

_____EXERCISES

**Exercise 19.1:**
    (i) The expected space for a skip list on $n$ keys is $\mathcal{O}(n)$.
    (ii) Describe an algorithm for insertion and analyze its expected cost.
    (iii) Describe an algorithm for deletion and analyze its expected cost.                    ◇

**Exercise 19.2:** The above (abstract) description of skip lists allows the possibility that $L_{i+1} = L_i$. So, the height $m$ can be arbitrarily large.
    (i) Pugh suggests some á priori maximum bound on the height (say, $m \leq k \lg n$) for some small constant $k$. Discuss the modifications needed to the algorithms and analysis in this case.
    (ii) Consider another simple method to bound $m$: if $L_{i+1} = L_i$ should happen, we simply omit the last key in $L_i$ from $L_{i+1}$. This implies $m \leq n$. Re-work the above algorithms and complexity analysis for this approach.                    ◇

**Exercise 19.3:** Determine the exact formula for $\mathbb{E}[X_{n,\ell}]$.                    ◇

**Exercise 19.4:** We have described skip lists in which each key is promoted to the next level with probability $p = 1/2$. Give reasons for choosing some other $p \neq 1/2$.                    ◇

## §20. Treaps

In our treatment of treaps, we assume that each item is associated with a *priority* as well as the usual *key*. A *treap* $T$ on a set of items is a binary search tree with respect to the keys of items, and a max-heap with respect to the priorities of items. Recall that a tree is a max-heap with respect to the priorities if the priority at any node is greater than the priorities in its descendents.

Suppose the set of keys and set of priorities are both $[1..n]$. Let the priority of key $k$ be denoted $\sigma(k)$. We assume $\sigma$ is a permutation of $[1..n]$. We explicitly describe $\sigma$ by listing the keys, in order of decreasing priority:

$$\sigma = (\sigma^{-1}(n), \sigma^{-1}(n-1), \ldots, \sigma^{-1}(1)). \tag{54}$$

We sometimes call $\sigma$ a *list* in reference to this representation. In general, if the range of $\sigma$ is not in $[1..n]$, we may define analogues of (54): a list of all the keys is called a $\sigma$-**list** if the priorities of the keys in the list are non-increasing. So $\sigma$-lists are unique if $\sigma$ assigns unique keys.

**Example.** Let $\sigma = (5, 2, 1, 7, 3, 6, 4)$. So the key 5 has highest priority of 7 and key 4 has priority 1. The corresponding treap is given by figure 6 where the key values are written in the nodes.
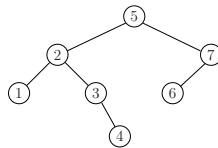


Figure 6: Treap $T_7(\sigma)$ where $\sigma = (5, 2, 1, 7, 3, 6, 4)$.

■

There is no particular relation between the key and the priority of an item. Informally, a *random treap* is a treap whose items are assigned priorities in some random manner.

FACT 1. *For any set $S$ of items, there is a treap whose node set is $S$. If the keys and priorities are unique for each item, the treap is unique.*

*Proof.* The proof is constructive: given $S$, pick the item $u_0$ with highest priority to be the root of a treap. The remaining items are put into two subsets that make up the left and right subtrees, respectively. The construction proceeds inductively. If keys and priorities are unique then $u_0$ and the partition into two subsets are unique.                    **Q.E.D.**

**Left paths and left spines.** The *left path* at a node $u_0$ is the path $(u_0, \ldots, u_m)$ where each $u_{i+1}$ is the left child of $u_i$ and the last node $u_m$ has no left child. The node $u_m$ is called the *tip* of the left path. The *left spine* at a node $u_0$ is the path $(u_0, u_1, \ldots, u_m)$ where $(u_1, \ldots, u_m)$ is the right path of the left child of $u_0$. If $u_0$ has no left child, the left spine is the trivial path $(u_0)$ of length 0. The *right path* and *right spine* is similarly defined. The *spine height* of $u$ is the sum of the lengths of the left and right spines of $u$. As usual, by identifying the root of a tree with the tree, we may speak of the left path, right spine, etc, of a tree. Note that the smallest item in a binary search tree is at the tip of the left path.

EXAMPLE: In figure 6, the left and right paths at node 2 are $(2, 1)$ and $(2, 3, 4)$, respectively. The left spine of the tree $(5, 2, 3, 4)$. The spine height of the tree is $3 + 2 = 5$.

## §21. Almost-treap and Treapification

We define a binary search tree $T$ to be an *almost-treap at $u$* $(u \in T)$ if $T$ is not a treap, but we can adjust the priority of $u$ so that $T$ becomes a treap. Roughly speaking, this means the heap property is satisfied everywhere except at $u$. Thus if we start with a treap $T$, we can change the priority of any node $u$ to get an almost-treap at $u$. We call an almost-treap an *under-treap* if the priority of $u$ must be adjusted upwards to get a treap, and an *over-treap* otherwise.

**Defects.** Suppose $T$ is an over-treap at $u$. Then there is a unique maximal prefix $\pi$ of the path from $u$ to the root such that the priority of each node in $\pi$ is less than the priority of $u$. The number of nodes in $\pi$ is called the *defect* in $T$. Note that $u$ is not counted in this prefix.

Next suppose $T$ is an under-treap at $u$. Then there are unique maximal prefixes of the left and right spines at $u$ such that the priority of each node in these prefixes is greater than the priority of $u$. The total number of nodes in these two prefixes is called the *defect* in $T$.

A pair $(u, v)$ of nodes is a *violation* if either $u$ is an ancestor of $v$ but $u$ has lower priority than $v$, or $u$ is a descendent of $v$ but with higher priority. It is not hard to check that that number of violations in an over-treap is equal to the defect. But this is not true in the case of an under-treap at $u$: this is because we only count violations along the two spines of $u$. We make some simple observations.

FACT 2. *Let $T$ be an almost-treap at $u$ with $k \geq 1$ defects. Let $v$ be a child of $u$ with priority at least that of its sibling (if any).*
    *(a) If $T$ is an over-treap then rotating at $u$ produces[10] an over-treap with $k - 1$ defects.*
    *(b) If $T$ is an over-treap, $v$ is defined and and $v$.priority $<$ $u$.priority then* rotate($v$) *produces an over-treap with $k + 1$ defects.*
    *(c) If $T$ be an under-treap and $v$ is defined then a rotation at $v$ results in an under-treap at $u$ with $k - 1$ defects.*

We now give a simple algorithm to convert an almost-treap into a treap. The argument to the algorithm is the node $u$ if the tree is an almost-treap at $u$. It is easy to determine from $u$ whether the tree is a treap or an under-treap or an over-treap.

---

[10] If $k = 1$, the result is actually a treap. But we shall accept the terminology "almost-treap with 0 defects" as a designation for a treap.

Algorithm Treapify$(u, T)$:
Input: $T$ is an almost-treap at $u$.
Output: a treap $T$.
　　　If $u$.priority $>$ $u$.parent.priority
　　　　　then Violation=OverTreap else Violation=UnderTreap.
　　　switch(Violation)
　　　　　case OverTreap:
　　　　　　　Do rotate$(u)$ until
　　　　　　　　　$u$.priority $\leq$ $u$.parent.priority.
　　　　　case UnderTreap:
　　　　　　　Let $v$ be the child of $u$ with the largest priority.
　　　　　　　While $u$.priority $<$ $v$.priority DO
　　　　　　　　　rotate$(v)$.
　　　　　　　　　$v \leftarrow$ the child of $u$ with largest priority.

We let the reader verify the following observation:

Fact 3. *Let $T$ be an almost-treap at $u$.*
*a) (Correctness) Treapify(u) converts $T$ into a treap.*
*b) If $T$ is an under-treap, then number of rotations is at most the spine height of $u$.*
*c) If $T$ is an over-treap, the number of rotations is at most the depth of $u$.*

――――――――――――――――――――――Exercises

**Exercise 21.1:** Let $n = 5$ and $\sigma = (3, 1, 5, 2, 4)$ *i.e.*, $\sigma(3) = 5$ and $\sigma(4) = 1$. Draw the treap. Next, change the priority of key 4 to 10 (from 1) and treapify. Next, change the priority of key 3 to 0 (from 5) and treapify. ◇

**Exercise 21.2:** Start with a treap $T$. Compare treapifying at $u$ after we increase the priority of $u$ by two different amounts. Is it true that the treapifying work is of the smaller increase is no more than the work for the larger increase? What if we decrease the priority of $u$ instead? ◇

## §22. Operations on Treaps

We show the remarkable fact that all the usual operations on binary search trees can be unified within a simple framework based on treapification.

In particular, we show how implement the following binary search tree operations:

　　　(a)　Lookup(Key,Tree)→Item,
　　　(b)　Insert(Item,Tree),
　　　(c)　Delete(Node,Tree),
　　　(d)　Successor(Item,Tree)→Item,
　　　(e)　Min(Tree)→Item,
　　　(f)　DeleteMin(Tree)→Item,
　　　(g)　Split(Tree1,Key)→Tree2,
　　　(h)　Join(Tree1,Tree2).

The meaning of these operations are defined as in the case of binary search trees. So the only twist is to simultaneously maintain the properties of a max-heap with respect to the priorities. First note that the usual binary tree operations of Lookup(Key,Tree), Successor(Item,Tree), Min(Tree) and Max(Tree) can be extended to treaps without change since these do not depend modify the tree structure. Next consider insertion and deletion.

(i) Insert$(u, T)$: we first insert in $T$ the item $u$ ignoring its priority. The result is an almost-treap at $u$. We now treapify $T$ at $u$.

(ii) Delete$(v, T)$: assuming no item in $T$ has priority $-\infty$, we first change the priority of the item at node $v$ into $-\infty$. Then we treapify the resulting under-treap at $v$. As a result, $v$ is now a leaf which we can delete directly.

We can implement DeleteMin(Tree) using Min(Tree) and Delete(item,Tree). It is interesting to observe that this deletion algorithm for treaps translates into a new deletion algorithm for ordinary binary search trees: to delete node $u$, just treapify at $u$ by pretending that the tree is an under-treap at $u$ until $u$ has at most one child. Since there are no priorities in ordinary binary search trees, we can rotate at any child $v$ of $u$ while treapifying.

**Split and join of treaps.**    Let $T$ be a treap on a set $S$ of items. We can implement these two operations easily using insert and deletes.

(i) *To split $T$ at $k$:* We insert a new item $u = (k, \infty)$ (*i.e.*, item $u$ has key $k$ and infinite priority). This yields a new treap with root $u$, and we can return the left subtree and right subtree to be the desired $T_L, T_R$.

(ii) *To join $T_L, T_R$:* first perform a Min$(T_R)$ to obtain the minimum key $k^*$ in $T_R$. Form the almost-treap whose root is the artificial node $(k^*, -\infty)$ and whose left and right subtrees are $T_L$ and $T_R$. Finally, treapify at $(k^*, -\infty)$ and delete it.

The beauty of treap algorithms is that they are reduced to two simple subroutines: binary insertion and treapification.

—————————————————————————————————————Exercises

**Exercise 22.1:** For any binary search tree $T$ and node $u \in T$, we can assign priorities to items such that $T$ is a treap and a deletion at $u$ takes a number of rotation equal to the the spine height of $u$.      ◇

**Exercise 22.2:** Modify the definition of the split operation so that all the keys in $S_R$ is strictly greater than the key $k$. Show how to implement this version.      ◇

**Exercise 22.3:** Define "almost-treap at $U$" where $U$ is a set of nodes. Try to generalize all the preceding results.      ◇

**Exercise 22.4:** Alternative deletion algorithm: above, we delete a node $u$ by making it a leaf. Describe an alternative algorithm where we step the rotations as soon as $u$ has only one

child. We then delete $u$ by replacing $u$ by its in-order successor or predecessor in the tree. Discuss the pros and cons of this approach.      ◇

## §23. Searching in a Random Treap

We analyze the expected cost of searching for a key in a random treap. Let us set up the random model. Assume the set of keys in the treap is $[1..n]$ and $S_n$ the set of permutations on $[1..n]$. The event space is $(S_n, 2^{S_n})$ with a uniform probability function $\Pr_n$. Each permutation $\sigma \in S_n$ represents a choice of priority for the keys $[1..n]$ where $\sigma(k)$ is the priority for key $k$. Recall our "list" convention (54) (§1) for writing $\sigma$.

The *random treap* on the keys $[1..n]$ is $T_n$ defined as follows: for $\sigma \in S_n$, $T_n(\sigma)$ is the treap on keys $[1..n]$ where the priorities of the keys are specified by $\sigma$ as described above. There are three ways in which such a sample space arise.
• We randomly generate $\sigma$ directly in a uniform manner. In section §6 we describe how to do this.
• We have a fixed but otherwise arbitrary continuous distribution function $F : \mathbb{R} \to [0,1]$ and the priority of key $i$ is a r.v. $X_i \in [0,1]$ with distribution $F$. The set of r.v.'s $X_1, \ldots, X_n$ are i.i.d.. In the continuous case, the probability that $X_i = X_j$ is negligible. In practice, we can approximate the range of $F$ by a suitably large finite set of values to achieve the same effect.
• The r.v.'s $X_1, \ldots, X_n$ could be given a uniform probability distribution if successive bits of each $X_i$ is obtained by tossing a fair coin. Moreover, we will generate as many bits of $X_i$ as are needed by our algorithm when comparing priorities. It can be shown that the expected number of bits needed is a small constant.

We frequently relate a r.v. of $\Pr_n$ to another r.v. of $\Pr_k$ for $k \neq n$. The following simple lemma illustrates a typical setting.

LEMMA 13. *Let $X, Y$ be r.v.'s of $\Pr_n, \Pr_k$ (respectively) where $1 \leq k < n$. Suppose there is a map $\mu : S_n \to S_k$ such that*
*(i) for each $\sigma' \in S_k$, the set $\mu^{-1}(\sigma')$ of inverse images has size $n!/k!$, and*
*(ii) for each $\sigma \in S_n$, $X(\sigma) = Y(\mu(\sigma))$.*
*Then $\mathrm{E}[X] = \mathrm{E}[Y]$.*

*Proof.*

$$\mathrm{E}[X] = \frac{\sum_{\sigma \in S_n} X(\sigma)}{n!} = \frac{\sum_{\sigma' \in S_k} (n!/k!) Y(\sigma')}{n!} = \frac{\sum_{\sigma' \in S_k} Y(\sigma')}{k!} = \mathrm{E}[Y].$$

     **Q.E.D.**

The map $\mu$ that "erases" from the string $\sigma$ all keys greater than $k$ has the properties stated in the lemma. E.g., if $n = 5$, $k = 3$ then $\sigma(3, 1, 4, 2, 5) = (3, 1, 2)$ and $\sigma^{-1}(3, 1, 2) = 4 \cdot 5 = 20$.

Recalling our example in figure 6: notice that if we insert the keys $[1..7]$ into an initially empty binary search tree in order of decreasing priority, (i.e., first insert key 5, then key 2, then key 1, etc), the result is the desired treap. This illustrates the following lemma:

LEMMA 14. *Suppose we repeatedly insert into a binary search tree the following sequence of keys $\sigma^{-1}(n), \sigma^{-1}(n-1), \ldots, \sigma^{-1}(1)$, in this order. If the initial binary tree is empty, the final binary tree is in fact the treap $T_n(\sigma)$.*

*Proof.* The proposed sequence of insertions results in a binary search tree, by definition. We only have to check that the heap property. This is seen by induction: the first insertion clearly results in a heap. If the $i$th insertion resulted in a heap, then the $(i+1)$st insertion is an almost heap. Since this newly inserted node is a leaf, and its rank is less than all the other nodes already in the tree, it cannot cause any violation of the heap property. So the almost heap is really a heap. Since treaps are unique for unique keys and ranks, this tree must be equal to $T_\sigma$.
**Q.E.D.**

The above lemma suggests that a random treap can be regarded as a binary search tree that arises from random insertions.

**Random ancestor sets $A_k$.** We are going to fix $k \in [1..n]$ as the key we are searching for in the random treap $T_n$. To analyze the expected cost of this search, we define the random set $A_k : S_n \to 2^{[1..n]}$ where

$$A_k(\sigma) := \{j : j \text{ is an ancestor of } k \text{ in } T_n(\sigma)\}.$$

By definition, $k \in A_k(\sigma)$. Clearly, the cost of $\text{LookUp}(T_n, k)$ is equal to $\Theta(|A_k|)$. Hence our goal is to determine $\mathbf{E}[|A_k|]$. To do this, we split the random variable $|A_k|$ into two parts:

$$|A_k| = |A_k \cap [1..k]| + |A_k \cap [k..n]| - 1.$$

By the linearity of expectation, we can determine the expectations of the two parts separately.

**Running $k$-maximas of $\sigma$.** We give a descriptive name to elements of the set

$$A_k(\sigma) \cap [1..k].$$

A key $j$ is called a *running $k$-maxima of $\sigma$* if $j \in [1..k]$ and for all $i \in [1..k]$,

$$\sigma(i) > \sigma(j) \Rightarrow i < j.$$

Of course, $\sigma(i) > \sigma(j)$ just means that $i$ appears before $j$ in the list

$$(\sigma^{-1}(n), \ldots, \sigma^{-1}(1)).$$

In other words, as we run down this list, by the time we come to $j$, it must be bigger than any other elements from the set $[1..k]$ that we have seen so far. Hence we call $j$ a "running maxima".

EXAMPLE. Taking $\sigma = (4, 3, 1, 6, 2, 7, 5)$ as in figure 6, the running 7-maximas of $\sigma$ are 4, 6 and 7. Note that no running $k$-maximas appears after $k$ in the list $\sigma$.

LEMMA 15. $j \in A_k(\sigma) \cap [1..k]$ *iff $j$ is a running $k$-maxima of $\sigma$.*

*Proof.* ($\Rightarrow$) Suppose $j \in A_k(\sigma) \cap [1..k]$. We want to show that $j$ is a running $k$-maxima. That is, for all $i \in [1..k]$ and $\sigma(i) > \sigma(j)$, we must show that $i < j$. Look at the path $\pi_j$ from the root to $j$ that is traced while inserting $j$. It will initially retrace the corresponding path $\pi_i$ for $i$ (which was inserted earlier). Let $i'$ be the last node that is common to $\pi_i$ and $\pi_j$ (so $i'$ is the least common ancestor of $i$ and $j$). We allow the possibility $i = i'$ (but surely $j \neq i'$ since $j = i'$ would make $j$ an ancestor of $i$, violating the max-heap property). There are two cases: $j$ is either (a) in the right subtree of $i'$ or (b) in the left subtree of $i'$. We claim that it cannot

be (b). To see this, consider the path $\pi_k$ traced by inserting $k$. Since $j$ is an ancestor of $k$, $\pi_j$ must be a prefix of $\pi_k$. If $j$ is a left descendent of $i'$, then so is $k$. Since $i$ is either $i'$ or lies in the left subtree of $i'$, this means $i > k$, contradiction. Hence (a) holds and we have $j > i' \geq i$.

($\Leftarrow$) Conversely, suppose $j$ is a running $k$-maxima in $\sigma$. It suffices to show that $j \in A_k(\sigma)$. This is equivalent to showing

$$A_j(\sigma) \subseteq A_k(\sigma).$$

View $T_n(\sigma)$ as the result of inserting a sequence of keys by decreasing priority. At the moment of inserting $k$, key $j$ has already been inserted. We just have to make sure that the search algorithm for $k$ follows the path $\pi_j$ from the root to $j$. This is seen inductively: clearly $k$ visits the root. Inductively, suppose $k$ is visiting a key $i$ on the path $\pi_j$. If $i > k$ then surely both $j$ and $k$ next visit the left child of $i$. If $i < k$ then $j > i$ (since $j$ is a running maxima) and hence again both $j$ and $k$ next visit the right child of $i$. This proves that $k$ eventually visits $j$.
                                                                                                    **Q.E.D.**

**Running $k$-minimas of $\sigma$.** Define a key $j$ to be a *running $k$-minima of $\sigma$* if $j \in [k..n]$ and for all $i \in [k..n]$,

$$\sigma(i) > \sigma(j) \Rightarrow j < i.$$

With $\sigma = (4, 3, 1, 6, 2, 7, 5)$ as before, the running 2-minimas of $\sigma$ are 4, 3 and 2. As for running $k$-maximas, no running $k$-minimas appear after $k$ in the list $\sigma$. We similarly have:

Lemma 16. $A_k(\sigma) \cap [k..n]$ *is the set of running $k$-minimas of $\sigma$.*

Define the random variable $U_n$ where $U_n(\sigma)$ is the number of running $n$-maximas in $\sigma \in S_n$.

We will show that $U_k$ has the same expected value as $|A_k \cap [1..k]|$. Note that $U_k$ is a r.v. of $\Pr_k$ while $|A_k \cap [1..k]|$ is a r.v. of $\Pr_n$. Similarly, define the r.v. $V_n$ where $V_n(\sigma)$ is the number of running 1-minimas in $\sigma \in S_n$.

Lemma 17.
*(a)* $\mathbb{E}[|A_k \cap [1..k]|] = \mathbb{E}[U_k]$.
*(b)* $\mathbb{E}[|A_k \cap [k..n]|] = \mathbb{E}[V_{n-k+1}]$.

*Proof.* Note that the r.v. $|A_k \cap [1..k]|$ is related to the r.v. $U_k$ exactly as described in lemma 13. This is because in each $\sigma \in S_n$, the keys in $[k+1..n]$ are clearly irrelevant to the number we are counting. Hence (a) is a direct application of lemma 13. A similar remark applies to (b) because the keys in $[1..k-1]$ are irrelevant in the running $k$-minimas.                    **Q.E.D.**

We give a recurrence for the expected number of running $k$-maximas.

Lemma 18.
*(a)* $\mathbb{E}[U_1] = 1$ *and* $\mathbb{E}[U_k] = \mathbb{E}[U_{k-1}] + \frac{1}{k}$ *for* $k \geq 2$. *Hence* $\mathbb{E}[U_k] = H_k$
*(b)* $\mathbb{E}[V_1] = 1$ *and* $\mathbb{E}[V_k] = \mathbb{E}[V_{k-1}] + \frac{1}{k}$ *for* $k \geq 2$. *Hence* $\mathbb{E}[V_k] = H_k$

*Proof.* We consider (a) only, since (b) is similar. Consider the map taking $\sigma \in S_k$ to $\sigma' \in S_{k-1}$ where we delete 1 from the sequence $(\sigma^{-1}(k), \ldots, \sigma^{-1}(1))$ and replace each remaining key $i$ by $i-1$, and take this sequence as the permutation $\sigma'$. For instance, $\sigma = (4, 3, 1, 5, 2)$ becomes

$\sigma' = (3, 2, 4, 1)$. Note that $U_5(\sigma) = 2 = U_4(\sigma')$. On the other hand, if $\sigma = (1, 4, 3, 5, 2)$ then $\sigma' = (3, 2, 4, 1)$. and $U_5(\sigma) = 3$ and $U_4(\sigma') = 2$. In general, we have

$$U_k(\sigma) = U_{k-1}(\sigma') + \delta(\sigma)$$

where $\delta(\sigma) = 1$ or $0$ depending on whether or not $\sigma(1) = k$. For each $\sigma' \in S_{k-1}$, that there are exactly $k$ permutations $\sigma \in S_k$ that maps to $\sigma'$; moreover, of these $k$ permutations, exactly one has $\delta(\sigma) = 1$. Thus

$$
\begin{aligned}
\mathrm{E}[U_k] &= \frac{\sum_{\sigma \in S_k} U_k(\sigma)}{k!} \\
&= \frac{\sum_{\sigma \in S_k} (U_{k-1}(\sigma') + \delta(\sigma))}{k!} \\
&= \frac{\sum_{\sigma' \in S_{k-1}} (1 + kU_{k-1}(\sigma'))}{k!} \\
&= \frac{1}{k} + \frac{\sum_{\sigma' \in S_{k-1}} U_{k-1}(\sigma')}{(k-1)!} \\
&= \frac{1}{k} + \mathrm{E}[U_{k-1}].
\end{aligned}
$$

Clearly, the solution to $\mathrm{E}[U_k]$ is the $k$th harmonic number $H_k = \sum_{i=1}^{k} 1/i$.      **Q.E.D.**

The depth of $k$ in the treap $T_n(\sigma)$ is

$$|A_k(\sigma)| - 1 = |A_k(\sigma) \cap [1..k]| + |A_k(\sigma) \cap [k..n]| - 2.$$

Since the cost of LookUp($k$) is proportional to 1 plus the depth of $k$, we obtain:

**Corollary 19.** *The expected cost of LookUp(k) in a random treap of n elements is proportional to*

$$\mathrm{E}[U_k] + \mathrm{E}[V_{n-k+1}] - 1 = H_k + H_{n-k+1} - 1 \leq 1 + 2\ln n.$$

_____EXERCISES

**Exercise 23.1:** (i) Prove lemma 16. (ii) Minimize $H_k + H_{n-k+1}$ for $k$ ranging over $[1..n]$.    ◇

**Exercise 23.2:**
    (i) Show that the variance of $U_n$ is $\mathcal{O}(\log n)$. In fact, $\mathrm{Var}(U_n) < H_n$.
    (ii) Use Chebyshev's inequality to show that the probability that $U_n$ is more than twice its expected value is $\mathcal{O}(\frac{1}{\log n})$.    ◇

## §24. Insertion and Deletion

A treap insertion has two phases:
*Insertion Phase.* This puts the item in a leaf position, resulting in an almost-treap.
*Rotation Phase.* This performs a sequence of rotations to bring the item to its final position. The insertion phase has the same cost as searching for the successor or predecessor of the item to be inserted. By the previous section, this work is expected to be $\Theta(\log n)$. What about the rotation phase?

Lemma 20. *For any set $S$ of items and $u \in S$, there is an almost-treap $T$ at $u$ in which $u$ is a leaf and the set of nodes is $S$. Moreover, if the keys and priorities in $S$ are unique, then $T$ is unique.*

*Proof.* To see the existence of $T$, we first construct the treap $T'$ on $S - \{u\}$. Then we insert $u$ into $T'$ to get an almost-treap at $u$. For uniqueness, we note that when the keys and priorities are unique and $u$ is a leaf of an almost-treap $T$, then the deletion of $u$ from $T$ gives a unique treap $T'$. On the other hand, $T$ is uniquely determined from $u$ and $T'$.      **Q.E.D.**

This lemma shows that insertion and deletion are inverses in a very strong sense. Recall that for deletion, we assume that we are given a node $u$ in the treap to be deleted. Then there are two parts again:

*Rotation Phase.* This performs a sequence of rotations to bring the node to a leaf position.
*Deletion Step.* This simply deletes the leaf.

We now clarify the precise sense in which rotation and deletion are inverses of each other.

**Corollary 21.** *Let $T$ be a treap containing a node $u$ and $T'$ be the treap after deleting $u$ from $T$. Let $T_+$ be the almost-treap just before the rotation phase while inserting $u$ into $T'$. Let $T_-$ be the almost-treap just after the rotation phase while deleting $u$ from $T$.*
    *(i) $T_+$ and $T_-$ are identical.*
    *(ii) The intermediate almost-treaps obtained in the deletion rotations are the same as the intermediate almost-treaps obtained in the insertion rotations (but in the opposite order).*

In particular, the costs of the rotation phase in deletion and in insertion are equal. Hence it suffices to analyze the cost of rotations in when deleting a key $k$. This cost is $\mathcal{O}(L_k + R_k)$ where $L_k$ and $R_k$ are the lengths of the left and right spines of $k$. These lengths are at most the depth of the tree, and we already proved that the expected depth is $\mathcal{O}(\log n)$. Hence we conclude: *the expected time to insert into or delete from a random treap is $\mathcal{O}(\log n)$.*

In this section, we give a sharper bound: the expected number of rotations during an insertion or deletion is at most 2.

**Random left-spine sets $B_k$.** Let us fix the key $k \in [1..n]$ to be deleted. We want to determine the expected spine height of $k$. We compute the left and right spine heights separately. For any $\sigma \in S_n$ and any key $k \in [1..n]$, define the set

$$B_k(\sigma) := \{j : j \text{ is in the left spine of } k\}.$$

**Triggered running $k$-maxima.** Again, we can give a descriptive name for elements in the random set $B_k$. Define $j$ to be a *triggered running $k$-maxima of $\sigma$* if $j \in [1..k-1]$ and $\sigma(k) > \sigma(j)$ and for all $i \in [1..k-1]$,

$$\sigma(i) > \sigma(j) \Rightarrow i < j.$$

In other words, as we run down the list

$$(\sigma^{-1}(n), \dots, k, \dots, j, \dots, \sigma^{-1}(1)),$$

by the time we come to $j$, we must have seen $k$ (this is the trigger) and $j$ must be bigger than any other elements in $[1..k-1]$ that we have seen so far. Note that the other elements in $[1..k-1]$ may occur before $k$ in the list (*i.e.*, they need not be triggered).

With $\sigma = (4, 3, 1, 6, 2, 7, 5)$ as in figure 6, the triggered running 3-maximas of $\sigma$ are 1 and 2. The triggered running 6-maximas of $\sigma$ is comprised of just 7; note that 2 is not a 6-maxima because 2 is less than some keys in $[1..5]$ which occurred earlier than 2.

LEMMA 22. *(Characterization of left-spine) Then $j \in B_k(\sigma)$ iff $j$ is a triggered running $k$-maxima of $\sigma$.*

*Proof.*($\Rightarrow$) Suppose $j \in B_k(\sigma)$. We want to show that
(i) $j < k$,
(ii) $\sigma(k) > \sigma(j)$, and
(iii) for all $i \in [1..k-1]$, $\sigma(i) > \sigma(j)$ implies $i < j$.
But (i) holds because of the binary search tree property, and (ii) holds by the heap property. What about (iii)? Let $i \in [1..k-1]$ and $\sigma(i) > \sigma(j)$. We want to prove that $i < j$. *First assume $i$ lies in the path from root to $j$.* There are two cases.
  CASE (1): $k$ is a descendent of $i$. Then we see that $k > i$ and $j$ is a descendent of $k$ implies $j > i$.
  CASE (2): $i$ is a descendent of $k$. But in a left-spine, the items have increasing keys as they lie further from the root. In particular, $i < j$, as desired.
*Now consider the case where $i$ does not lie in the path from the root to $j$.* Let $i'$ be the least common ancestor of $i$ and $j$. Clearly either

$$i < i' < j \qquad \text{or} \qquad j < i' < i$$

must hold. So it suffices to prove that $i' < j$ (and so the first case hold). Note that $i' \neq k$. If $i'$ is a descendent of $k$ then essentially CASE (2) above applies, and we are done. So assume $k$ is a descendent of $i'$. Since $j$ is also a descendent of $k$, it follows that $i$ and $k$ lie in different (left or right) subtrees of $i'$. By definition $i < k$. Hence $i$ lies in the left subtree of $i'$ and $k$ lies in the right subtree of $i'$. Hence $j$ lies in the right subtree of $i'$, i.e., $j > i'$, as desired. This proves (iii).
($\Leftarrow$) Suppose $j$ is a triggered running $k$-maxima. We want to show that $j \in B_k(\sigma)$. We view $T_n(\sigma)$ as the result of inserting items according to priority (§3): assume that we are about to insert key $j$. Note that $k$ is already in the tree. We must show that $j$ ends up in the "tip" of the left-spine of $k$. For each key $i$ along the path from the root to the tip of the left-spine of $k$, we consider two cases.
  CASE (3): $i > k$. Then $i$ is an ancestor of $k$ and (by (i)) we have $i > k > j$. This means $j$ will move into the subtree of $i$ which contains $k$.
  CASE (4): $i < k$. Then by (iii), $i < j$. Again, this means that if $i$ is an ancestor of $k$, $j$ will move into the subtree of $i$ that contains $k$, and otherwise $i$ is in the left-spine of $k$ and $j$ will move into the right subtree of $i$ (and thus remain in the left-spine of $k$).      **Q.E.D.**

It is evident that the keys in $[k+1..n]$ is irrelevant to the number of triggered running $k$-maximas. Hence we may as well assume that $\sigma \in S_k$ (instead of $S_n$). To capture this, let us define the r.v. $T_k$ that counts the number of triggered running $k$-maximas in case $n = k$.

LEMMA 23. *Consider the r.v. $|B_k|$ of $\Pr_n$. It is related to the r.v. $T_k$ of $\Pr_k$ via $\mathbb{E}[|B_k|] = \mathbb{E}[T_k]$.*

*Proof.*For each $\sigma \in S_n$, let us define its *image* $\sigma'$ *in* $S_k$ obtained by taking the sequence $(\sigma^{-1}(n), \ldots, \sigma^{-1}(1))$, deleting all keys of $[k+1..n]$ from this sequence, and considering the result as a permutation $\sigma' \in S_k$. Note that $T_k(\sigma') = |B_k(\sigma)|$. It is not hard to see that there are exactly $n!/k!$ choices of $\sigma \in S_k$ whose images equal a given $\sigma' \in S_k$. Hence lemma 13 implies $\mathbb{E}[|B_k|] = \mathbb{E}[T_k]$.      **Q.E.D.**

LEMMA 24. *We have $\mathbb{E}[T_1] = 0$. For $k \geq 2$,*

$$\mathbb{E}[T_k] = \mathbb{E}[T_{k-1}] + \frac{1}{k(k-1)}$$

*and hence* $\mathbf{E}[T_k] = \frac{k-1}{k}$.

*Proof.*It is immediate that $\mathbf{E}[T_1] = 0$. So assume $k \geq 2$. Again consider the map from $\sigma \in S_k$ to a corresponding $\sigma' \in S_{k-1}$ where $\sigma'$ is the sequence obtained from $(\sigma^{-1}(k), \ldots, \sigma^{-1}(1))$ by deleting the key 1 and subtracting 1 from each of the remaining keys. Note that

- each $\sigma' \in S_{k-1}$ is the image of exactly $k$ distinct $\sigma \in S_k$, and
- for each $j \in [2..k-1]$, $j$ is a triggered running $k$-maxima in $\sigma$ iff $j-1$ is a triggered running $(k-1)$-maxima in $\sigma'$.
- key 1 is a triggered running $k$-maxima in $\sigma$ iff $\sigma(k) = k$ and $\sigma(1) = k - 1$.

This proves
$$T_k(\sigma) = T_{k-1}(\sigma') + \delta(\sigma)$$

where $\delta(\sigma) = 1$ or 0 depending on whether or not $\sigma = (k, 1, \sigma^{-1}(k-2), \ldots, \sigma^{-1}(1))$. Now, there are exactly $(k-2)!$ choices of $\sigma \in S_k$ such that $\delta(\sigma) = 1$. Hence

$$
\begin{aligned}
\mathbf{E}[T_k] &= \frac{\sum_{\sigma \in S_k} T_k(\sigma)}{k!} \\
&= \frac{\sum_{\sigma \in S_k} (T_{k-1}(\sigma') + \delta(\sigma))}{k!} \\
&= \frac{\sum_{\sigma \in S_k} T_{k-1}(\sigma')}{k!} + \frac{(k-2)!}{k!} \\
&= \frac{\sum_{\sigma' \in S_{k-1}} kT_{k-1}(\sigma')}{k!} + \frac{1}{k(k-2)} \\
&= \mathbf{E}[T_{k-1}] + \frac{1}{k(k-2)}.
\end{aligned}
$$

It is immediate that
$$\mathbf{E}[T_k] = \sum_{i=2}^{k} \frac{1}{i(i-1)},$$

using the fact that $\mathbf{E}[T_1] = 0$. This is a telescoping sum in disguise because

$$\frac{1}{i(i-1)} = \frac{1}{i-1} - \frac{1}{i}.$$

The solution follows at once.        **Q.E.D.**

**Right spine.** Using symmetry, the expected length for the right spine of $k$ is easily obtained. We define $j$ to be a *triggered running $k$-minimas in $\sigma$* if (i) $j > k$, (ii) $\sigma(j) < \sigma(k)$, (iii) for all $i \in [k+1..n]$, $\sigma(i) > \sigma(j)$ implies $j < i$. We leave as an exercise to show:

Lemma 25. *$j$ is in the right spine of $k$ iff $j$ is a triggered running $k$-minima of $\sigma$.*

The keys in $[1..k-1]$ are irrelevant for this counting. Hence we may define the random variable
$$R_k(\sigma)$$

which counts the number of triggered running 1-minimas in $\sigma \in S_k$. We have:

Lemma 26.
(a) *The expected length of the right spine is given by $\mathbf{E}[R_{n-k+1}]$.*
(b) *For $k \geq 2$, $\mathbf{E}[R_k] = \mathbf{E}[R_{k-1}] + \frac{1}{k(k-1)}$. The solution is $\mathbf{E}[R_k] = (k-1)/k$.*

---

**Corollary 27.** *The expected length of the right spine of $k$ is*

$$\mathrm{E}[R_{n-k+1}] = \frac{n-k}{n-k+1}.$$

Hence the expected lengths of the left and right spines of a key $k$ is less than 1 each. This leads to the surprising result:

**Corollary 28.** *The expected number of rotations in deleting or inserting a key from the random treap $T_n$ is less than 2.*

This is true because each key is likely to be a leaf or near one. So this result is perhaps not surprising since more than half of the keys are leaves.

_____Exercises

**Exercise 24.1:** Show lemma 25.      ◇

## §25. Cost of a Sequence of Requests

We have shown that for single requests, it takes expected $\mathcal{O}(\log n)$ time to search or insert, and $\mathcal{O}(1)$ time to delete. Does this bound apply to a sequence of such operations? The answer is not obvious because it is unclear whether the operations may upset our randomness model. How do we ensure that inserts or deletes to a random treap yields another random treap?

We formalize the problem as follows: fix any sequence

$$p_1, \ldots, p_m$$

of treap requests on the set $[1..n]$ of keys. In other words, each $p_i$ ($k \in [1..n]$ has one of the forms

$$\text{LookUp}(k), \text{Insert}(k), \text{Delete}(k).$$

We emphasize that that the above sequence of requests is arbitrary, not "random" in any probabilistic sense. Also note that we assume that deletion operation is based on a key value, not on a "node" as we normally postulate (§4). This variation should not be a problem in most applications.

Our probability space is again

$$(S_n, 2^{S_n}, \mathrm{Pr}_n).$$

For each $i = 1, \ldots, m$ and $\sigma \in S_n$, let $T_i(\sigma)$ be the treap obtained after the sequence of operations $p_1, \ldots, p_i$, where $\sigma$ specifies the priorities assigned to keys and $T_0(\sigma)$ is the initial empty treap. Define $X_i$ to be the r.v. such that $X_i(\sigma)$ is the cost of performing the request $p_i$ on the treap $T_{i-1}(\sigma)$, resulting in $T_i(\sigma)$. The expected cost of operation $p_i$ is simply $\mathrm{E}[X_i]$.

Lemma 29.
$$\mathrm{E}[X_i] = \mathcal{O}(\log i).$$

*Proof.*Fix $i = 1, \ldots, m$. Let $K_i \subseteq [1..n]$ denote the set of keys that are in $T_i(\sigma)$. A key observation is the dependence of $T_{i-1}(\sigma)$ on $p_1, \ldots, p_{i-1}$ amounts only to the fact that these $i-1$ operations determine $K_i$. Let $|K_i| = n_i$. Then there is a natural map taking $\sigma \in S_n$ to $\sigma' \in S_{n_i}$ (*i.e.*, the map deletes from the sequence $(\sigma^{-1}(n), \ldots, \sigma^{-1}(1))$ all elements not in $K_i$, and finally replacing each remaining key $k \in K_i$ by a corresponding key $\pi_i(k) \in [1..n_i]$ which preserves key-order). Each $\sigma'$ is the image of $n!/(n_i)!$ distinct $\sigma \in S_n$. Moreover, for any key $k \in K_i$, its relevant statistics (the depth of $k$ and length of left/right spines of $k$) in $T_i(\sigma)$ is the same as that of its replacement key $\pi_i(k)$ in the treap on $[1..k]$. As shown above, the expected values of these statistics is the same as the expected values of these statistics on the uniform probability space on $S_{n_i}$. But we have proven that the expected depth of any $k \in [1..n_i]$ is $\mathcal{O}(\log n_i)$, and the expected length of the spines of $k$ is at most 1. This proves that when the update operation is based on a key $k \in K_i$, $\mathrm{E}[X_i] = \mathcal{O}(\log n_i) = \mathcal{O}(\log i)$, since $n_i \leq i$.

What if the operation $p_i$ is an unsuccessful search, i.e., we do a LookUp$(k)$ where $k \notin K_i$? In this case, the length of the path that the algorithm traverses is bounded by the depth of either the successor or predecessor of $k$ in $K_i$. Again, the expected lengths in $\mathcal{O}(\log i)$. This concludes our proof.        **Q.E.D.**

**Application to Sorting.** We give a randomized sorting algorithm as follows: on input a set of $n$ keys, we assign them a random priority (by putting them in an array in decreasing order of priority). Then we do a sequence of $n$ inserts in order of this priority. Finally we do a sequence of $n$ DeleteMins. This gives an expected $\mathcal{O}(n \log n)$ algorithm for sorting.

NOTES: Galli and Simon (1996) suggested an alternative approach to treaps which avoid the use of random priorities. Instead, they choose the root of the binary search tree randomly. Mulmuley introduced several abstract probabilistic "game models" that corresponds to our analysis of running (triggered) maximas or minimas.

—————————————————————————Exercises

**Exercise 25.1:** What is the worst case and best case key for searching and insertion in our model?        ◇

—————————————————————————End Exercises

# References

[1] N. Alon, J. H. Spencer, and P. Erdös. *The probabilistic method.* John Wiley and Sons, Inc, New York, 1992.

[2] C. R. Aragon and R. G. Seidel. Randomized search trees. *IEEE Foundations of Comp. Sci.*, 30:540–545, 1989.

[3] H. Chernoff. A measure of asymptotic efficiency for tests of hypothesis based on sum of observations. *Ann. of Math. Stat.*, 23:493–507, 1952.

[4] K. L. Chung. *Elementary Probability Theory with Stochastic Processes.* Springer-Verlag, New York, 1979.

[5] W. Feller. *An introduction to Probability Theory and its Applications.* Wiley, New York, 2nd edition edition, 1957. (Volumes 1 and 2).

[6] A. T. Jonassen and D. E. Knuth. A trivial algorithm whose analysis isn't. *J. Computer and System Sciences*, 16:301–322, 1978.

[7] R. M. Karp. Probabilistic recurrence relations. *J. ACM*, 41(6):1136–1150, 1994.

[8] D. E. Knuth. *The Art of Computer Programming: Fundamental Algorithms*, volume 1. Addison-Wesley, Boston, 2nd edition edition, 1975.

[9] D. E. Knuth. *The Art of Computer Programming: Seminumerical Algorithms*, volume 2. Addison-Wesley, Boston, 2nd edition edition, 1981.

[10] A. N. Kolmogorov. *Foundations of the theory of probability.* Chelsea Publishing Co., New York, 1956. Second English Edition.

[11] M. Li and P. M. B. Vitányi. *An introduction to Kolmogorov Complexity and its Applications.* Springer-Verlag, second edition, 1997.

[12] R. Motwani and P. Raghavan. *Randomized Algorithms.* Cambridge University Press, 1995.

[13] W. W. Peterson and J. E. J. Weldon. *Error-Correcting Codes.* MIT Press, 1975. 2nd Edition.

[14] W. Pugh. Skip lists: a probabilistic alternative to balanced trees. *Commun. ACM*, 33(6):668–676, 1990.