

# GAMBLING TEAMS AND WAITING TIMES FOR PATTERNS IN TWO-STATE MARKOV CHAINS

JOSEPH GLAZ, MARTIN KULLDORFF,

VLADIMIR POZDNYAKOV, AND J. MICHAEL STEELE

ABSTRACT. Methods using gambling teams and martingales are developed and applied to find formulas for the expected value and the generating function of the waiting time until one observes an element of a finite collection of patterns in a sequence which is generated by a two-state first or higher order Markov chain. (KEYWORDS: Gambling, teams, waiting times, patterns, success runs, failure runs, Markov chains, martingales, stopping times, generating functions. MATHEMATICS SUBJECT CLASSIFICATION (2000): Primary 60J10, Secondary 60G42)

## 1. INTRODUCTION

How long must one observe a stochastic process with values from a finite alphabet until one sees a realization of a pattern which belongs to a specified collection  $\mathcal{C}$  of possible patterns? For independent processes this is an old question; in some special cases it is even considered by Feller (1968). Nevertheless, in the context of more

---

*Date:* November 9, 2005.

J. Glaz and V. Pozdnyakov: Department of Statistics, University of Connecticut, 215 Glenbrook Road, U-4120, Storrs, CT 06269-4120.

M. Kulldorff: Department of Ambulatory Care and Prevention, Harvard Medical School and Harvard Pilgrim Health Care, 133 Brookline Avenue, Boston, MA 02215-3920.

J. M. Steele: Wharton School, Department of Statistics, Huntsman Hall 447, University of Pennsylvania, Philadelphia, PA 19104.

general process, or even for Markov chains, there are many natural problems which have not been fully addressed. The main goal here is to show how some progress can be made by further developing the martingale methods which were introduced by Li (1980) and Li and Gerber (1981) in their investigation of independent sequences. Their key observation was that information on the occurrence times of patterns can be obtained from the values assumed by a specially constructed auxiliary martingale at a certain well-chosen time.

In the case of (first- or higher-order) Markov chains, this observation is still useful, but to make it work requires a rather more elaborate plan for the construction of the auxiliary martingale. This construction depends in turn on several general devices which seem more broadly useful; these include “teams of gamblers,” “watching then betting,” “reward matching,” and a couple of other devices which will be described shortly.

Before engaging that description, we should note that pattern matching has been studied by many other techniques. The combinatorial methods of Guibas and Odlyzko (1981a, 1981b) are particularly effective, and there are numerous treatments of pattern matching problem by probabilistic techniques, such as Benevento (1984), Biggins and Cannings (1987a, 1987b), Blom and Thorburn (1982), Breen et al. (1985), Chrysaphinou and Papastavridis (1990), Han and Hirano (2003), Pozdnyakov et al. (2005), Pozdnyakov and Kulldorff (2006), Robin and Daudin (1999), Stefanov (2003) and Uchida (1998). One of the more general techniques is the Markov chain embedding method introduced by Fu (1986) which has been further developed by Antzoulakos (2001), Fu (2001), Fu and Chang (2002), and Fu and Koutras (1994). The approach of Stefanov (2000) and Stefanov and Pakes (1997)

also use Markov chain embedding, though their method differs substantially from Fu's. Only a few investigations considered waiting time problems for higher order Markov chains, and these have all focused on specific waiting times such as the "sooner or later" problem studied by Aki et al. (1996).

## 2. EXPECTED WAITING TIME UNTIL A PATTERN IS OBSERVED

We take  $\{Z_n, n \geq 1\}$  to be a Markov chain with two states  $S$  and  $F$ , which may model "success" and "failure." We suppose the chain has the initial distribution  $\mathbf{P}(Z_1 = S) = p_S$ ,  $\mathbf{P}(Z_1 = F) = p_F$  and the transition matrix

$$\begin{pmatrix} p_{SS} & p_{FS} \\ p_{SF} & p_{FF} \end{pmatrix},$$

where  $p_{SF}$  is shorthand for  $\mathbf{P}(Z_{n+1} = F | Z_n = S)$ . We then consider a collection  $\mathcal{C}$  of finite sequences  $A_i$ ,  $1 \leq i \leq K$ , from the two-letter alphabet  $\{S, F\}$ . If  $\tau_{A_i}$  denotes the first time until the pattern  $A_i$  has been observed as a completed run in the series  $Z_1, Z_2, \dots$ , then the random variable of main interest here is  $\tau_{\mathcal{C}} = \min\{\tau_{A_1}, \dots, \tau_{A_K}\}$ , the first time when we observe a pattern from  $\mathcal{C}$ . Throughout our discussion we assume that no pattern of  $\mathcal{C}$  contains another pattern from  $\mathcal{C}$  as an initial segment. Naturally, this assumption entails no loss of generality.

**2.1. A Run of "Failures" Under a Markov Model.** To illustrate the construction, we first consider the rather easy case where  $K = 1$  and were the pattern  $A_1$  is a run of  $r$  consecutive  $F$ s. Thus, we will compute the expected value of  $\tau = \tau_{A_1}$ , the time of the first completion of a run of  $r$  "failures" under our two-state Markov model. This example can be handled by several methods, and it offers a useful benchmark for more challenging examples.

We consider a casino where gamblers may bet in successive rounds on the output of our given two-state Markov chain, and we assume that the casino is *fair* in a sense which we will soon make precise. We then consider a sequence of gamblers, one of whom arrives just before each new round of betting. Thus, gambler number  $n + 1$  arrives in time to observe the result of the  $n^{\text{th}}$  trial,  $Z_n$ , and we assume that he bets a dollar on the event that next trial yields an  $F$ . If  $Z_{n+1} = S$ , he loses his dollar and leaves the game. If he is lucky and  $Z_{n+1} = F$ , then he wins  $1/p_{SF}$  when  $Z_n = S$  and he wins  $1/p_{FF}$  when  $Z_n = F$ . This is the sense in which the casino is fair; the expected return on a one dollar bet is one dollar.

After this gambler gets his money, he then bets his entire capital on the event that  $Z_{n+2} = F$ . Again, if  $Z_{n+2} = S$ , then the gambler leaves the game with nothing. On the other hand, if  $Z_{n+2} = F$ , then the gambler wins this round, and his capital is increased by the factor  $1/p_{FF}$ . Successive rounds proceed in the same way, with a new gambler arriving at each new round and with the gamblers from earlier periods either continuing to win or else going broke and leaving.

Now we need to be precise about the end of this process. If gambler  $n + 1$  begins by observing  $Z_n = S$ , then he bets until either he goes broke or until he observes  $r$  successive  $F$ s, and, if gambler  $n + 1$  begins by observing  $Z_n = F$ , then bets until he either goes broke or until he observes  $r - 1$  successive  $F$ s. Once some gambler stops without going broke, all of the gambling stops.

Finally, we let  $X_n$  denote the casino's net gain at conclusion of round  $n$ . Since each bet is fair and since the bet sizes depend only on the previous observations, the sequence  $X_n$  is a martingale with respect to the  $\sigma$ -field generated by  $\{Z_n, n \geq 1\}$ . Now we just need to consider the casino's net gain  $X_\tau$  when the gambling stops.

By calculating  $E(X_\tau)$  in two ways we will then obtain the expected value of the time  $\tau$ .

At time  $\tau$  many gamblers are likely to have lost all their money; only those who entered the game after round number  $\tau - r - 1$  have any money. We now face two different ending scenarios. First, it could happen that we have a block (denoted by  $F^{(r)}$ ) of  $r$  instances of  $F$  which occur at the very beginning of the sequence  $\{Z_n, n \geq 1\}$ . Second, it could happen that we end with  $SF^{(r)}$ , an  $S$  followed by a block of  $r$  instances of  $F$ . Obviously we do not need to consider the possibility of ending with  $FF^{(r)} = F^{(r+1)}$  since by definition  $F^{(r)}$  cannot occur before time  $\tau$ .

When we total up the wins and losses of all of the gamblers, we then find that the value of the stopped martingale  $X_\tau$  is given exactly by

$$X_\tau = \begin{cases} \tau - 1 - \frac{1}{p_{FF}^{r-1}} - \frac{1}{p_{FF}^{r-2}} - \dots - \frac{1}{p_{FF}}, & \text{1st scenario,} \\ \tau - 1 - \frac{1}{p_{SF}p_{FF}^{r-1}} - \frac{1}{p_{FF}^{r-1}} - \frac{1}{p_{FF}^{r-2}} - \dots - \frac{1}{p_{FF}}, & \text{2nd scenario,} \end{cases}$$

which can be written more briefly as

$$X_\tau = \begin{cases} \tau - 1 - \frac{1 - p_{FF}^{r-1}}{p_{FF}^{r-1}(1 - p_{FF})}, & \text{1st scenario,} \\ \tau - 1 - \frac{1}{p_{SF}p_{FF}^{r-1}} - \frac{1 - p_{FF}^{r-1}}{p_{FF}^{r-1}(1 - p_{FF})}, & \text{2nd scenario.} \end{cases}$$

Since  $\mathbf{E}[\tau] < \infty$  and the increments of  $X_n$  are bounded, the optional stopping theorem for martingales (for instance, Williams (1991, p. 100)) tells us that  $0 = \mathbf{E}[X_1] = \mathbf{E}[X_\tau]$ . From this identity and the formula for  $X_\tau$ , algebraic simplification gives us

$$(1) \quad \mathbf{E}[\tau] = 1 + p_F \frac{1 - p_{FF}^{r-1}}{(1 - p_{FF})} + (1 - p_F p_{FF}^{r-1}) \left( \frac{1}{p_{SF}p_{FF}^{r-1}} + \frac{1 - p_{FF}^{r-1}}{p_{FF}^{r-1}(1 - p_{FF})} \right).$$

**2.2. Second Step: A Single Pattern.** We now consider the more subtle case of a single (non-run) pattern  $A$  with length  $r$ , and for specificity we assume that the pattern begins with  $F$ , so  $A = FB$  where we have  $B \in \{S, F\}^{r-1}$ . As before we consider a sequence of gamblers, but this time we need to consider three different *ending scenarios*:

- (1)  $A$  occurs at the beginning of the sequence  $\{Z_n, n \geq 1\}$ , or
- (2) the pattern  $SA$  occurs, or
- (3) the pattern  $FA$  occurs.

The probability  $p_1$  of the first scenario is trivial to compute, but one then runs into trouble. We do not know the probability that the pattern  $SA$  will appear earlier than  $FA$ , so the probabilities of the second and third ending scenarios are not readily computed. To circumvent this problem we introduce two teams of gamblers.

**2.3. Rules for the Gambling Teams.**

- (1) A gambler from the first team who arrives before round  $n$  watches the result of the  $n^{\text{th}}$  trial, and then bets  $y_1$  dollars on the first letter in the sequence  $A$ . If he wins he then bets all of his capital on the next letter in the sequence  $A$ , and he continues in this way until he either loses his capital or he observes all of the letters of  $A$ . Such players are called *straightforward gamblers*.
- (2) The gamblers of the second team make use of the information that they observe. If gambler  $n + 1$  observes  $Z_n = S$  just before he begins his play, then he bets just like a straightforward gambler except that he begins by wagering  $y_2$  dollars on the first letter of pattern  $A$ . On the other hand, if he observes  $Z_n = F$  when he first arrives, then wagers  $y_2$  dollars on the

first letter of the pattern  $B$ . He then continues to wager on the successive letters of  $B$  either until he loses or until he observes  $B$ . Such players are called *smart gamblers*.

The two gambling teams continue their betting, until one team stops. At that time, all gambling stops, and we consider the wins and losses. Only those gamblers who enter the game after the time  $\tau - r - 1$  will have any money and the amount they have will depend on the ending scenario. If we let  $W_{ij}y_j$  denote the amount of money that team  $j \in \{1, 2\}$  wins in scenario  $i \in \{1, 2, 3\}$ , then the values  $W_{ij}$  are easy to compute, and in terms of these values of stopped martingale  $X_\tau$  which represents the casino's net gain is given by

$$X_\tau = \begin{cases} (y_1 + y_2)(\tau - 1) - y_1W_{11} - y_2W_{12}, & 1^{\text{st}} \text{ scenario,} \\ (y_1 + y_2)(\tau - 1) - y_1W_{21} - y_2W_{22}, & 2^{\text{nd}} \text{ scenario,} \\ (y_1 + y_2)(\tau - 1) - y_1W_{31} - y_2W_{32}, & 3^{\text{rd}} \text{ scenario.} \end{cases}$$

Now, if we take  $(y_1^*, y_2^*)$  to be a solution of the system

$$y_1^*W_{21} + y_2^*W_{22} = 1, \quad y_1^*W_{31} + y_2^*W_{32} = 1,$$

we see that with these bet sizes we have a very simple formula for  $X_\tau$ :

$$X_\tau = \begin{cases} (y_1^* + y_2^*)(\tau - 1) - y_1^*W_{11} - y_2^*W_{12}, & 1^{\text{st}} \text{ scenario,} \\ (y_1^* + y_2^*)(\tau - 1) - 1, & 2^{\text{nd}} \text{ scenario,} \\ (y_1^* + y_2^*)(\tau - 1) - 1, & 3^{\text{rd}} \text{ scenario.} \end{cases}$$

The optional stopping theorem then gives us

$$0 = (y_1^* + y_2^*)(\mathbf{E}[\tau] - 1) - p_1(y_1^*W_{11} + y_2^*W_{12}) - (1 - p_1),$$

where  $p_1$  is the probability of scenario one. We therefore find

$$(2) \quad \mathbf{E}[\tau] = 1 + \frac{p_1(y_1^*W_{11} + y_2^*W_{12}) + (1 - p_1)}{y_1^* + y_2^*}.$$

Formula (2) is more explicit than it may seem at first. In the typical case, the calculation of  $p_1$ ,  $\{W_{ij} : 1 \leq i \leq 3, 1 \leq j \leq 2\}$  and  $\{y_j^* : 1 \leq j \leq 2\}$  is genuinely routine, as one can see by the next example.

**2.4. Example: Waiting Time Until FSF.** Here our straightforward gamblers bet  $y_1$  dollars on  $FSF$  without regard of the preceding observation. On the other hand, the smart gamblers bet  $y_2$  dollars on  $FSF$  if they observed  $S$  before placing their first bet, but they bet  $y_2$  dollars on  $SF$  if they observed  $F$ . The three ending scenarios are now either  $FSF$  at the beginning (scenario one), or one ends with  $SFSF$  (scenario two), or one ends with  $FFSF$  (scenario three). The  $3 \times 2$  “profit matrix”  $\{W_{ij}\}$  is then given by

$$\begin{pmatrix} \frac{1}{p_{SF}} & \frac{1}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \\ \frac{1}{p_{SF}p_{FS}p_{SF}} + \frac{1}{p_{SF}} & \frac{1}{p_{SF}p_{FS}p_{SF}} + \frac{1}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \\ \frac{1}{p_{FF}p_{FS}p_{SF}} + \frac{1}{p_{SF}} & \frac{1}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \end{pmatrix},$$

and bet sizes  $y_1^*$  and  $y_2^*$  are determined by the relations

$$\begin{aligned} y_1^* \left( \frac{1}{p_{SF}p_{FS}p_{SF}} + \frac{1}{p_{SF}} \right) + y_2^* \left( \frac{1}{p_{SF}p_{FS}p_{SF}} + \frac{1}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \right) &= 1, \\ y_1^* \left( \frac{1}{p_{FF}p_{FS}p_{SF}} + \frac{1}{p_{SF}} \right) + y_2^* \left( \frac{1}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \right) &= 1, \end{aligned}$$

which one can solve to obtain

$$y_1^* = \frac{p_{FF}p_{FS}p_{SF}}{p_{FS} + p_{SF} + p_{FS}p_{SF}} \quad \text{and} \quad y_2^* = \frac{p_{FS}p_{SF}(p_{SF} - p_{FF})}{p_{FS} + p_{SF} + p_{FS}p_{SF}}.$$



The probability  $p_1$  of the first scenario is just  $p_F p_{FS} p_{SF}$ , so after substitution and simplification the general formula (2) provides

$$\mathbf{E}[\tau_{FSF}] = 1 + \frac{p_S}{p_{SF}} + \frac{1}{p_{SF}^2} + \frac{1}{p_{FS} p_{SF}},$$

which is as explicit as one could wish.

### 3. EXPECTED TIME UNTIL OBSERVING ONE OF MANY PATTERNS

We now consider a collection  $\mathcal{C} = \{A_i : 1 \leq i \leq K\}$  of  $K$  strings of possibly varying lengths from the two-letter alphabet, and we take on the task of computing the expected value of  $\tau_{\mathcal{C}} = \min\{\tau_{A_1}, \dots, \tau_{A_K}\}$ , the first time that one observes one of the patterns in  $\mathcal{C}$ . The method we propose is analogous to the two-team method we just used, although many teams are now needed. The real challenge is the construction of the list of the appropriate ending scenarios which now requires some algorithmic considerations.

**3.1. Listing the Ending Scenarios.** Given  $\mathcal{C} = \{A_i\}_{1 \leq i \leq K}$  we first consider the set sequence transformation

$$\mathcal{C} = \{A_i\}_{1 \leq i \leq K} \longrightarrow \{SA_i, FA_i\}_{1 \leq i \leq K} = \{B_i\}_{1 \leq i \leq 2K} = \mathcal{C}',$$

which doubles the cardinality of  $\mathcal{C}$ . We then delete from  $\mathcal{C}'$  each pattern  $B$  which can only occur after the stopping time  $\tau_{\mathcal{C}}$ . The resulting collection  $\mathcal{C}''$  is called the *final list*. We denote the elements of  $\mathcal{C}''$  by  $C_i$ ,  $1 \leq i \leq K'$ , and we note that  $K \leq K' \leq 2K$ .

To illustrate the construction, suppose the initial collection is  $\mathcal{C} = \{FSF, FF\}$ . The *doubling step* gives us  $\mathcal{C}' = \{SFSF, FFSS, SFF, FFF\}$ . Since  $FFS$  and  $FFF$  cannot occur before  $\tau$ , these are eliminated from  $\mathcal{C}'$  and the final list is

simply  $\mathcal{C}'' = \{SF SF, SFF\}$ . Similarly, if the initial collection is  $\mathcal{C} = \{FS, SSS\}$ , then the final list is  $\mathcal{C}'' = \{SFS, FFS\}$ .

Now, before we describe the ending scenarios, we need one further notion. If patterns  $C$  and  $C'$  in the final list  $\mathcal{C}''$  satisfy  $C = SA$  and  $C' = FA$  for some pattern  $A \in \mathcal{C}$ , then we say that  $C$  and  $C'$  are *matched*. Also, if  $C$  and  $C'$  are matched and  $C = SA$  and  $C' = FA$ , then we say that  $C$  and  $C'$  are *generated* by  $A$ . Finally, even though there are many ending scenarios, they may be classified to three basic kinds.

- (1) There are  $K$  scenarios where one observes an element of  $\mathcal{C}$  as an initial segment of the Markov sequence  $\{Z_n, n \geq 1\}$ .
- (2) There is a scenario for each unmatched pattern from  $\mathcal{C}''$ . We denote the number of these by  $L$ .
- (3) There is pair of scenarios for each matched pattern from  $\mathcal{C}''$ . We denote the number of these by  $2M$ .

**3.2. From the Listing to the Teams.** For each scenario associated with unmatched pattern  $C_j$  we introduce *one* team of straightforward gamblers who bet  $y_j$  dollars on the pattern  $A_i$  which generated  $C_j$ . For each pair of scenarios associated with matched patterns  $C_p$  and  $C_m$  which were generated by pattern  $A_k$ , we introduce *two* teams. One team bets  $y_p$  dollars on  $A_k$  in the *straightforward* way, another bets  $y_m$  dollars on  $A_k$  in the *smart* way. If  $W_{ij}y_j$ ,  $i = 1, 2, \dots, K + L + 2M, j = 1, 2, \dots, L + 2M$  denotes amount of money that the  $j^{\text{th}}$  team wins in the  $i^{\text{th}}$  scenario, then the stopped martingale  $X_\tau$  is given by the sum

$$X_\tau = \sum_{j=1}^{L+2M} y_j(\tau - 1) - S(y_1, \dots, y_{L+2M}),$$

where we set

$$S(y_1, \dots, y_{L+2M}) = \sum_{i=1}^{K+L+2M} 1_{E_i} \sum_{j=1}^{L+2M} y_j W_{ij},$$

and where  $1_{E_i}$  is the indicator function for the event  $E_i$  that the  $i^{\text{th}}$  scenario occurs.

If  $(y_1^*, \dots, y_{L+2M}^*)$  is a solution of the linear system

$$(3) \quad \begin{aligned} y_1^* W_{K+1 \ 1} + \dots + y_{L+2M}^* W_{K+1 \ L+2M} &= 1, \\ \vdots & \\ y_1^* W_{K+L+2M \ 1} + \dots + y_{L+2M}^* W_{K+L+2M \ L+2M} &= 1, \end{aligned}$$

then we have

$$S(y_1^*, \dots, y_{L+2M}^*) = \begin{cases} \sum_{j=1}^{L+2M} y_j^* W_{ij}, & \text{in scenario } i \in \{1, 2, \dots, K\} \\ 1, & \text{in scenario } i > K \end{cases}$$

By the optional stopping theorem we have

$$0 = \mathbf{E}[X_1] = \mathbf{E}[X_{\tau_C}] = \sum_{j=1}^{L+2M} y_j^* (\mathbf{E}[\tau_C] - 1) - \sum_{i=1}^K p_i \sum_{j=1}^{L+2M} y_j^* W_{ij} - (1 - \sum_{i=1}^K p_i),$$

where  $p_i$  is the probability that  $A_i$  is an initial segment of  $\{Z_n, n \geq 1\}$ . We can now solve this equation to obtain a formula for  $\mathbf{E}[\tau_C]$  which we summarize as a theorem.

**Theorem 1.** *If  $(y_1^*, y_2^*, \dots, y_{L+2M}^*)$  solves the linear system (3), then*

$$(4) \quad \mathbf{E}[\tau_C] = 1 + \frac{\sum_{i=1}^K p_i \sum_{j=1}^{L+2M} y_j^* W_{ij} + (1 - \sum_{i=1}^K p_i)}{\sum_{j=1}^{L+2M} y_j^*}.$$

As before, this formula is more explicit than it may seem at first. In particular example, all of the required terms can be computed straightforwardly in problems of reasonable size.

**3.3. Computation of the Profit Matrix.** Formula (4) requires one to compute the profit matrix  $\{W_{ij}\}$ , and we will first show how this can be done in general. The method will then be applied to a specific example to confirm that formula (4) may be rewritten in terms of the basic model parameters.

Consider a scenario that ends with pattern  $C = c_1c_2\dots c_m \in \mathcal{C}''$ . The team of straightforward gamblers who begin by betting one dollar and who bet on the successive terms of the pattern  $A = a_1a_2\dots a_p$  will by time  $\tau_C$  have won

$$\sum_{i=1}^{\min(m-1,p)} \delta_i^{st}(A, C),$$

where

$$\delta_i^{st}(A, C) = \begin{cases} \frac{1}{p_{c_{m-i}a_1}p_{a_1a_2}\dots p_{a_{i-1}a_i}}, & \text{if } a_1 = c_{m-i+1}, a_2 = c_{m-i+2}, \dots, a_i = c_m \\ 0, & \text{otherwise.} \end{cases}$$

Similarly, the team of smart gamblers will have won

$$\sum_{i=1}^{\min(m-1,p)} \delta_i^{sm1}(A, C) + \sum_{i=1}^{\min(m-1,p-1)} \delta_i^{sm2}(A, C),$$

where we set

$$\delta_i^{sm1}(A, C) = \begin{cases} \frac{1}{p_{c_{m-i}a_1}p_{a_1a_2}\dots p_{a_{i-1}a_i}}, & \text{if } a_1 = c_{m-i+1}, a_2 = c_{m-i+2}, \dots, \\ & a_i = c_m \text{ and } c_{m-i} \neq a_1 \\ 0, & \text{otherwise,} \end{cases}$$

and

$$\delta_i^{sm2}(A, C) = \begin{cases} \frac{1}{p_{a_1a_2}p_{a_2a_3}\dots p_{a_i a_{i+1}}}, & \text{if } a_1 = c_{m-i}, a_2 = c_{m-i+1}, \dots, a_{i+1} = c_m \\ 0, & \text{otherwise.} \end{cases}$$

**3.4. Explicit Determination of  $\mathbf{E}[\tau_{\mathcal{C}}]$ .** To illustrate the use of formula (4), we consider  $\mathcal{C} = \{SS, FSF\}$ . After doubling and elimination we get the final list  $\{FSS, SFSS, FFSF\}$ , and we then need to work out the set of scenarios. We have two scenarios where  $C_1 = SS$  or  $C_2 = FSF$  occur as an initial segment of  $\{Z_n, n \geq 1\}$ . We also have the unmatched scenario  $C_3 = FSS$  associated with the pattern  $SS$ , and we have a pair of matched scenarios  $C_4 = SFSS$  or  $C_5 = FFSF$  which are associated with the pattern  $FSF$ . The profit matrix  $\{W_{ij}\}$  is then given by

$$\begin{pmatrix} \frac{1}{p_{SS}} & 0 & 0 \\ 0 & \frac{1}{p_{SF}} & \frac{1}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \\ \frac{1}{p_{FS}p_{SS}} + \frac{1}{p_{SS}} & 0 & 0 \\ 0 & \frac{1}{p_{SF}p_{FS}p_{SF}} + \frac{1}{p_{SF}} & \frac{1}{p_{SF}p_{FS}p_{SF}} + \frac{1}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \\ 0 & \frac{1}{p_{FF}p_{FS}p_{SF}} + \frac{1}{p_{SF}} & \frac{1}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \end{pmatrix},$$

and, after solving the corresponding linear system, we find that the appropriate initial team bets are given by

$$y_1^* = \frac{p_{FS}p_{SS}}{1 + p_{FS}}, \quad y_2^* = \frac{p_{FF}p_{FS}p_{SF}}{p_{FS} + p_{SF} + p_{FS}p_{SF}}, \quad y_3^* = \frac{p_{FS}p_{SF}(p_{SF} - p_{FF})}{p_{FS} + p_{SF} + p_{FS}p_{SF}}.$$

The probabilities  $p_1$  and  $p_2$  that  $SS$  and  $FSF$  are initial segments of the process  $\{Z_n, n \geq 1\}$  are given by  $p_{SS}$  and  $p_{FF}p_{FS}p_{SF}$  respectively, so the formula (4) leads one to the pleasantly succinct result

$$\mathbf{E}[\tau_{\mathcal{C}}] = 2 + p_{SF} + \frac{1 - p_{SS}}{p_{FS}}.$$

## 4. GENERATING FUNCTIONS FOR PATTERN WAITING TIMES

To find the generating function of the waiting time  $\tau_{\mathcal{C}}$  we need to introduce the same scenarios and the same betting teams, but we need to make some changes in the design of the initial bets. A gambler from the  $j^{\text{th}}$  team who arrives at moment  $k-1$  and who begins his betting on round  $k$  will now begin with a bet of size  $y_j \alpha^k$  where  $0 < \alpha < 1$ . If  $\alpha^{\tau_{\mathcal{C}}} W_{ij}(\alpha) y_j$  denotes total winnings of the  $j^{\text{th}}$  team when the game ends with  $i^{\text{th}}$  scenario, then we call  $W_{ij}(\alpha)$  the  $\alpha$ -profit matrix. As before, the  $\alpha$ -profit matrix does not depend on  $\tau_{\mathcal{C}}$ , and it can be computed if we know the ending scenario.

If  $X_n$  again denotes the casino's net gain at moment  $n$ , then

$$X_{\tau_{\mathcal{C}}} = \frac{\alpha^2 - \alpha \alpha^{\tau_{\mathcal{C}}}}{1 - \alpha} \sum_{j=1}^{L+2M} y_j - S(\alpha, y_1, \dots, y_{L+2M}),$$

and we set

$$S(\alpha, y_1, \dots, y_{L+2M}) = \sum_{i=1}^{K+L+2M} 1_{E_i} \sum_{j=1}^{K+2M} \alpha^{\tau_{\mathcal{C}}} y_j W_{ij}(\alpha),$$

where, as before,  $1_{E_i}$  is the indicator function for the event  $E_i$  that the  $i^{\text{th}}$  scenario occurs.

If  $(y_1^*, \dots, y_{L+2M}^*)$  is a solution of the linear system

$$(5) \quad \begin{aligned} y_1^* W_{K+1 \ 1}(\alpha) + \dots + y_{L+2M}^* W_{K+1 \ L+2M}(\alpha) &= 1, \\ \vdots & \\ y_1^* W_{K+L+2M \ 1}(\alpha) + \dots + y_{L+2M}^* W_{K+L+2M \ L+2M}(\alpha) &= 1, \end{aligned}$$

then we might hope to mimic our earlier calculation of  $\mathbf{E}[X_{\tau_{\mathcal{C}}}]$ , but unfortunately we run into trouble since  $\mathbf{E}(\alpha^{\tau_{\mathcal{C}}} 1_{E_1})$  may not equal  $p_1 \mathbf{E} \alpha^{\tau_{\mathcal{C}}}$ .

Nevertheless, if the  $i^{\text{th}}$  with  $i \leq K$  scenario occurs, then we know exactly the value of  $\tau_{\mathcal{C}}$ . It is equal to  $|A_i|$  — the length of  $i^{\text{th}}$  sequence. Therefore, we have a

formula for the stopped martingale,

$$X_{\tau_c} = \frac{\alpha^2 - \alpha\alpha^\tau}{1 - \alpha} \sum_{j=1}^{L+2M} y_j^* - \alpha^\tau - I(\alpha, y_1^*, \dots, y_{L+2M}^*),$$

where  $I(\alpha, y_1^*, \dots, y_{L+2M}^*)$  is defined by setting

$$I(\alpha, y_1^*, \dots, y_{L+2M}^*) = \begin{cases} \alpha^{|A_i|} \left[ \sum_{j=1}^{L+2M} y_j^* W_{ij}(\alpha) - 1 \right], & \text{in scenario } i \in \{1, 2, \dots, K\} \\ 0, & \text{in scenario } i > K. \end{cases}$$

From this formula, the optional stopping theorem, we then find an the anticipated formula for the moment generating function of  $\tau_c$ .

**Theorem 2.** *If  $(y_1^*, \dots, y_{L+2M}^*)$  is a solution of linear system (5), then one has*

$$(6) \quad \mathbf{E}[\alpha^{\tau_c}] = \frac{\frac{\alpha^2}{1-\alpha} \sum_{j=1}^{L+2M} y_j^* - \sum_{i=1}^K p_i \alpha^{|A_i|} \left[ \sum_{j=1}^{L+2M} y_j^* W_{ij} - 1 \right]}{1 + \frac{\alpha}{1-\alpha} \sum_{j=i}^{L+2M} y_j^*}.$$

**4.1. Computation of the  $\alpha$ -Profit Matrix.** As before, one needs to know how to compute the profit matrix, before formula (6) may be properly regarded as an explicit formula. This is only a little more difficult than before. First, assume that a scenario ends with the pattern  $C = c_1 c_2 \dots c_m$ . The team of straightforward gamblers who bet a dollar on pattern  $A = a_1 a_2 \dots a_p$  by the time  $\tau$  will win

$$\sum_{i=1}^{\min(m-1, p)} \delta_i^{st}(A, C) / \alpha^{i-1},$$

while the team of smart gamblers will win

$$\sum_{i=1}^{\min(m-1, p)} \delta_i^{sm1}(A, C) / \alpha^{i-1} + \sum_{i=1}^{\min(m-1, p-1)} \delta_i^{sm2}(A, C) / \alpha^{i-1}.$$

These formulas provide almost everything we need, but before we can be completely explicit, we need to focus on a concrete example.

**4.2. A Generating Function Example.** Consider the waiting time until one observes the 3-letter pattern  $FSF$  in the random sequence  $\{Z_n, n \geq 1\}$  produced by the Markov model. In this case, the  $\alpha$ -profit matrix  $\{W_{ij}\}$  is given by

$$\begin{pmatrix} \frac{1}{p_{SF}} & \frac{\alpha^{-1}}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \\ \frac{\alpha^{-2}}{p_{SF}p_{FF}p_{SF}} + \frac{1}{p_{SF}} & \frac{\alpha^{-2}}{p_{SF}p_{FF}p_{SF}} + \frac{\alpha^{-1}}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \\ \frac{\alpha^{-2}}{p_{FF}p_{FF}p_{SF}} + \frac{1}{p_{SF}} & \frac{\alpha^{-1}}{p_{FS}p_{SF}} + \frac{1}{p_{SF}} \end{pmatrix},$$

and by solving the associated linear system one finds

$$y_1^* = \frac{\alpha^2 p_{FF}p_{FF}p_{SF}}{1 - \alpha p_{FF} + \alpha p_{SF} + \alpha^2 p_{FF}p_{SF}}, \quad y_2^* = \frac{\alpha^2 p_{FS}p_{SF}(p_{SF} - p_{FF})}{1 - \alpha p_{FF} + \alpha p_{SF} + \alpha^2 p_{FF}p_{SF}}.$$

The general moment generating representation (6) then gives us the simple formula

$$\mathbf{E}[\alpha^{\tau_C}] = \frac{\alpha^3 p_{FS}p_{SF}(p_F + \alpha(p_S - p_{SS}))}{1 - \alpha(p_{SS} + p_{FF} - \alpha(p_{FF} - p_{SF}(1 - p_{FS}(1 - \alpha p_{SS})))}.$$

Naturally, such a formula provides one with complete information on the distribution of  $\tau_C$ , and to obtain an explicit formula for  $P(\tau_C = k)$  one can use symbolic calculation to rewrite the rational function (6) in its partial fraction expansion.

## 5. HIGHER ORDER MARKOV CHAINS

Here we have applied the gambling team method only to two-state chains, and, for reasons which will be explained later, this limitation is not easily lifted. Nevertheless, there are more complex chains where the team method applies, and it is instructive to consider one of these. Specifically, we briefly consider how the gambling team method may be applied with second order two-state chains. Here we obviously need to avoid the naive representation of such chains as first order chains with four states.



In the team approach for a second order model the gamblers need to observe *two* rounds of betting before they place their first bets, and consequently we need to consider a larger number of final scenarios. Moreover, for each pattern  $A = a_1 a_2 \dots a_p$  we will need to consider up to seven termination cases, including three “initial” cases which are associated with the patterns (1)  $A$ , (2)  $SA$ , or (3)  $FA$  and four “later” cases which are associated with the patterns (4)  $SSA$  (5)  $SFA$ , (6)  $FSA$ , and (7)  $FFA$ .

As before our main objective is to count accurately all the ending scenarios and create a matched number of gambling teams. However, in the second order chain case there is an additional difficulty that one needs to address. More specifically, one needs to consider separately two cases: (1) there are no runs in the initial list  $\mathcal{C}$  and (2) there are runs in  $\mathcal{C}$ .

**5.1. The First Case: There are no Runs in  $\mathcal{C}$ .** First we need to replace the earlier doubling step with an analogous *quadrupling* step. Now given the collection  $\mathcal{C} = \{A_i\}_{1 \leq i \leq K}$  of patterns, we consider the set sequence transformation

$$\mathcal{C} = \{A_i\}_{1 \leq i \leq K} \longrightarrow \{SSA_i, SFA_i, FSA_i, FFA_i\}_{1 \leq i \leq K} = \{B_i\}_{1 \leq i \leq 4K} = \mathcal{C}'.$$

We then delete from  $\mathcal{C}'$  each scenario which can happen only after the stopping time  $\tau_{\mathcal{C}}$ , and we take the collection  $\mathcal{C}''$  that remains to be our “final list” of ending scenarios.

Each pattern from the collection  $\mathcal{C}$  leads us to four — or perhaps fewer — ending scenarios. Now for each pattern from  $\mathcal{C}$  we consider a sequence of gamblers who belong to teams of different types. As before, these gamblers arrive sequentially, and they observe the game before placing any bets.

- (1) A new gambler from the *type I* team arrives two rounds before he begins to bet. He watches these rounds and then bets on the successive letters pattern  $A$ , with complete indifference to what he may have seen on the first two rounds.
- (2) A new gambler from the *type II* team also arrives two rounds before he can begin to gamble, but he is influenced by what he sees. If this gambler observes  $Sa_1$  on these two rounds, then he bets on the sequence  $a_2a_3\dots a_p$ , but, if he observes anything other than  $Sa_1$ , then he places his bets according to the sequence  $A$ .
- (3) Similarly, a new gambler from the *type III* team watches two rounds, and if he observes  $Fa_1$  then he bets according to the sequence  $a_2a_3\dots a_p$ , but, if he observes anything other than  $Fa_1$ , then he places his bets according to  $A$ .
- (4) Finally, a gambler from the *type IV* team watches two rounds, and if he observes  $a_1a_2$  then he bets according to the sequence  $a_3a_4\dots a_p$ ; otherwise he bets according to the sequence  $A$ .

Each pattern  $A$  from initial list  $\mathcal{C}$  leads to zero, one, two, three or four scenarios in the final list  $\mathcal{C}''$ , so now instead of just having to consider matched and unmatched patterns the patterns in final list  $\mathcal{C}''$  are of four kinds: *unmatched*, *double-matched*, *triple-matched*, and *quadruple-matched*.

To see how this works, consider the initial collection  $\mathcal{C} = \{FSFF, FFSF\}$ . First, note that we have five initial cases, (1)  $FSFF$ , (2)  $FFSF$ , (3)  $SFSFF$ , (4)  $SFFSF$ , (5)  $FFFSF$ . Pattern  $FFFSF$  cannot occur before  $\tau_{\mathcal{C}}$ , therefore, the scenario associated with this pattern is eliminated from the list of initial cases. Second,

in the final list  $\mathcal{C}''$  we have five patterns: double-matched patterns  $SSFSSF$  and  $FSFSSF$  generated by  $FSFF$  and triple-matched patterns  $SSFFSF$ ,  $SFFFSF$  and  $FFFFSF$  generated by  $FFSF$ . Thus, we need to introduce five teams: type I and II teams that bet on  $FSFF$ , and type I, II and III teams that bet on  $FFSF$ .

**5.2. The Second Case: Special Treatment of Runs.** If the initial list  $\mathcal{C}$  contains a run then one may have a problem with straightforward application of the method described above. The difficulty is that if we observe the game only till moment  $\tau_{\mathcal{C}}$  then there is no difference in behavior between teams of different types that place their bets on the run.

To illustrate the problem let us consider the initial list  $\mathcal{C} = \{F^{(r)}\}$ . The straightforward usage of the above algorithm tells us that one has to introduce two initial cases, (1)  $F^{(r)}$  and (2)  $SF^{(r)}$ . The final list  $\mathcal{C}''$  contains two double-matched patterns  $SSF^{(r)}$  and  $FSF^{(r)}$ . Therefore, according to the algorithm one needs two (type I and II) teams that bet  $y_1$  and  $y_2$  dollars on  $F^{(r)}$ . However, since before time  $\tau_{\mathcal{C}}$  there is no difference in gambling between these two teams, one, in fact, has just one team that bets  $y_1 + y_2$  dollars on the run. Thus, the number of free parameters is not matching the number of ending scenarios.

But a simple modification of the gambling method easily solves the problem. Before time  $\tau_{\mathcal{C}}$  the run  $F^{(r)}$  can only occur as an initial segment of the sequence  $\{Z_n, n \geq 1\}$  or as pattern  $SF^{(r)}$  later. So, if the initial list  $\mathcal{C}$  contains runs (obviously we can have one or two runs in  $\mathcal{C}$  only –  $F^{(r)}$  or  $S^{(p)}$  or both), then we need first to substitute runs  $F^{(r)}$  and  $S^{(p)}$  in  $\mathcal{C}$  by  $SF^{(r)}$  and  $FS^{(p)}$ , respectively, to get a different collection  $\tilde{\mathcal{C}}$ . The collection  $\tilde{\mathcal{C}}$  contains no runs, therefore, we can proceed as before. After application of quadrupling and elimination processes to the list  $\tilde{\mathcal{C}}$

we will get the final list of ending scenarios  $\tilde{\mathcal{C}}''$ , and for this list we will be able to create a matched number of gambling teams. Since we are interested in  $\tau_{\mathcal{C}}$  not  $\tau_{\tilde{\mathcal{C}}}$ , the elimination process has to be based on  $\tau_{\mathcal{C}}$ , and the runs must be included in the list of initial cases.

For instance, if  $\mathcal{C} = \{F^{(r)}\}$  one needs to consider four initial cases (1)  $F^{(r)}$ , (2)  $SF^{(r)}$ , (3)  $SSF^{(r)}$  and (4)  $FSF^{(r)}$ , and four later cases where the game ends by (5)  $SSSF^{(r)}$ , (6)  $SFSF^{(r)}$ , (7)  $FSSF^{(r)}$ , or (8)  $FFSF^{(r)}$ . In this case one can show that all four teams (type I, II, III and IV) that bet on  $SF^{(r)}$  bet in its own way.

**5.3. Final Step.** After attending to this bookkeeping, we can now calculate the expected observation times in a way that parallels our earlier calculation. Since we have matched the number of (non-initial) ending scenarios and the number of teams, we can choose the size of initial bet for each team in a way that makes all the expressions for the stopped martingale equal to 1 — however the game may end.

Let us summarize this as a theorem. Assume that in the end we have  $P$  initial cases and  $Q$  later cases. Let  $W_{ij}y_j$ ,  $i = 1, 2, \dots, P+Q$ ,  $j = 1, 2, \dots, Q$  denotes amount of money that the  $j^{\text{th}}$  team that bets  $y_j$  dollars wins in the  $i^{\text{th}}$  scenario. Finally, let  $p_i$ ,  $i = 1, 2, \dots, P$  be the probability that the  $i^{\text{th}}$  initial case takes place.

**Theorem 3.** *If  $(y_1^*, y_2^*, \dots, y_Q^*)$  solves the linear system*

$$\begin{aligned} y_1^* W_{P+1\ 1} + \cdots + y_Q^* W_{P+1\ Q} &= 1, \\ \vdots & \\ y_1^* W_{P+Q\ 1} + \cdots + y_Q^* W_{P+Q\ Q} &= 1, \end{aligned}$$

then

$$\mathbf{E}[\tau_{\mathcal{C}}] = 2 + \frac{\sum_{i=1}^P p_i \sum_{j=1}^Q y_j^* W_{ij} + (1 - \sum_{i=1}^P p_i)}{\sum_{j=1}^Q y_j^*}.$$

## 6. CONCLUDING REMARK

The method of gambling teams deals quite effectively with the waiting time problems of two-state chains, but for  $N$ -state chains, it is much less effective. The problem is that typically one finds that the number of ending scenarios is higher than the number of teams one has, so there are too few free parameters to achieve the requested matching.

One might think of reducing the waiting time problems for an  $N$ -state chains by encoding the states  $\{1, 2, \dots, N\}$  as sequences of zeros and ones, but this idea typically fails since the natural encodings do not lead one to a waiting time problem for a *homogeneous* two-state Markov chain on  $\{0, 1\}$ . Ironically, for many of the pattern problems associated with  $N$ -state Markov chains, the method of gambling team is ineffective when  $N \geq 3$ , even though for the corresponding problems in a two-state chain, it is typically the method of choice.

A possible *computational* advantage of the martingale method over the Markov chain embedding method (e.g., Antzoulakos (2001), Fu (2001), Fu and Chang (2002)) is the size of matrices involved in the calculation. The size of the profit matrix depends only on the number of patterns  $K$ , while the size of the transition matrix of embedded Markov chain also depends on the length of patterns from  $\mathcal{C}$ . For example, if  $\mathcal{C}$  contains  $K$  patterns each of which has a length that is about  $N$ , and  $K$  is much smaller than  $N$ , then the dimension of the transition matrix in the Markov chain embedding method is about  $K \times N$  by  $K \times N$ . For large  $N$  this size can cause technical problems. The size of profit matrix is at most  $2K$  by  $2K$ .

## REFERENCES

- [1] Aki, S., Balakrishnan, N. and Mohanty, S.G. (1996). Sooner and later waiting time problems and failure runs in higher order Markov dependent trials, *Ann. Inst. Statist. Math.*, **48**, 773-787.
- [2] Antzoulakos, D. (2001). Waiting times for patterns in a sequence of multistate trials, *J. Appl. Prob.* **38**, 508-518.
- [3] Benevento, R.V. (1984). The occurrence of sequence patterns in ergodic Markov chains. *Stoc. Proc. Appl.*, **17**, 369-373.
- [4] Biggins, J.D. and Cannings C. (1987a). Formulas for mean restriction-fragment lengths and related quantities. *Am. J. Hum. Genet.*, **41**, 258-265.
- [5] Biggins, J. D. and Cannings, C. (1987b). Markov renewal processes, counters and repeated sequences in Markov chains, *Adv. Appl. Prob.*, **19**, 521-545.
- [6] Blom, G. and Thorburn, D. (1982). How many random digits are required until given sequences are obtained?, *J. Appl. Prob.*, **19**, 518-531.
- [7] Breen, S., Waterman, M., and Zhang, N. (1985). Renewal theory for several patterns, *J. Appl. Prob.*, **22**, 228-234.
- [8] Chrysaphinou, O. and Papastavridis, S. (1990). The occurrence of a sequence of patterns in repeated dependent experiments, *Theory of Probability and Applications*, **35**, 145-152.
- [9] Feller, W (1968). *An introduction to probability theory and its applications, Vol 1, 3rd ed.* Wiley, New York
- [10] Fu, J. C. (1986). Reliability of consecutive- $k$ -out-of- $n$ :  $F$  systems with  $(k - 1)$ -step Markov dependence, *IEEE Trans. Reliability*, **R35**, 602-606.
- [11] Fu, J. (2001). Distribution of the scan statistics for a sequence of bistate trials, *J. Appl. Prob.*, **38**, 908-916.
- [12] Fu, J. and Chang, Y. (2002). On probability generating functions for waiting time distribution of compound patterns in a sequence of multistate trials, *J. Appl. Prob.*, **39**, 70-80.
- [13] Fu, J. C. and Koutras, M. V. (1994) Distribution theory of runs: A Markov chain approach, *J. Amer. Statist. Assoc.* **78**, 168-175.
- [14] Guibas, L. and Odlyzko, A. (1981a). Periods of strings, *J. Comb. Theory, Ser. A*, **30**, 19-42.

- [15] Guibas, L. and Odlyzko, A. (1981b). String overlaps, pattern matching and nontmsitive games, *J. Comb. Theory. Ser. A*, **30**, 183-208.
- [16] Gerber, H. and Li, S. (1981). The occurrence of sequence patterns in repeated experiments and hitting times in a Markov chain, *Stoch. Proc. Appl.*, **11**, 101-108.
- [17] Han, Q. and Hirano, K. (2003). Sooner and later waiting time problems for patterns in Markov dependent trials, *J. Appl. Probab.*, **40**, 73-86.
- [18] Li, S. (1980). A martingale approach to the study of occurrence of sequence patterns in repeated experiments, *Ann. Prob.*, **8**, 1171-1176.
- [19] Pozdnyakov, V., Glaz, J., Kulldorff, M., and Steele, J. M. (2005). A martingale approach to scan statistics, *Ann. Inst. Statist. Math.*, **57**, 21-37.
- [20] Pozdnyakov, V., and Kulldorff, M. (2006). Waiting Times for Patterns and a Method of Gambling Teams, *The American Mathematical Monthly*, **113**, 134-143.
- [21] Robin, S. and Daudin, J.-J. (1999). Exact distribution of word occurrences in a random sequence of letters, *J. Appl. Probab.*, **36**, 179-193
- [22] Stefanov, V. T. (2000). On some waiting time problems, *J. Appl. Prob.*, **37**, 756-764.
- [23] Stefanov, V. T. (2003). The intersite distances between pattern occurrences in strings generated by general discrete- and continuous-time models: an algorithmic approach, *J. Appl. Prob.*, **40**, 881-892.
- [24] Stefanov, V. T. and Pakes, A. G. (1997). Explicit distributional results in pattern formation, *Ann. Appl. Prob.*, **7**, 666-678.
- [25] Uchida, M. (1998). On generating functions of waiting time problems for sequence patterns of discrete random variables, *Ann. Inst. Statist. Math.*, **50**, 655-671.
- [26] Williams, D. (1991). *Probability with martingales*, Cambridge University Press, Cambridge.