

Using a specialized Markov chain in the reliability model of disk arrays RAID-10 with data mirroring and striping

P A Rahman

Department of Automated Technological and Informational Systems, Ufa State Petroleum Technological University, 2, October ave., Sterlitamak, 453118, Russia

E-mail: pavelar@yandex.ru

Abstract. This paper deals with a specialized Markov chain for a reliability model of the fault-tolerant system with dual-type failures, non-synchronized repairs and full restore after reaching the failed state. A generalized iterative procedure for calculation of the system's stationary availability factor, the mean time of repair and the mean time of failure is introduced. An application of the specialized Markov chain in the reliability model of the fault-tolerant dual-level disk arrays 'RAID-10' and calculation of the mean time to data loss by using the generalized iterative procedure are also observed. Finally, calculation examples of mean time to data loss of the RAID-10 array are also given.

1. Introduction

Currently, the world deals with complex technical devices and systems, which are used in everyday life as well as in the industrial manufacturing processes. In addition to these parameters of technical systems, like performance, time latency, capacity, power consumption etc., the reliability parameters are also quite important, because they directly affect the stability, efficiency and security of technical systems.

There is a set of the academic books, dedicated to reliability theory [1-4], but these in most cases deal with the well-known reliability model of the repairable systems, based on Markov Birth-Death Chain [5, 6]. These models allow us to obtain calculation formulas for such complex and integral reliability parameters like mean time to failure and a stationary availability factor. However, in more sophisticated cases, with several types of failures and special kinds of repairs, well-known Markov Birth-Death Chain cannot be applied to the reliability model of these systems, and development of the specialized Markov Chains is required. In particular, the reliability model of the fault-tolerant dual-level disk array 'RAID-10' (RAID – redundant array of inexpensive disks) [7, 8] cannot be properly represented by Markov Birth-Death Chain and requires development of a specialized chain.

Over the past few years, within the ambit of the research work in the field of reliability models of data transmission, processing and storage systems [9, 10], the author developed a specialized Markov Chain for the reliability model of the fault-tolerant system with n elements, dual-type failures, non-synchronized repairs and full restore after reaching the failed state, and obtained a generalized procedure for calculation of the mean time to failure and a stationary availability factor.

Finally, the author applied the specialized Markov Chain to the reliability model of the fault-tolerant dual-level 'RAID-10' array and obtained a specialized procedure for calculation of the mean time to data loss of the 'RAID-10' array.



2. Specialized Markov chain for the reliability model of the system with dual-type failures, non-synchronized repairs and full restore after reaching the failed state

Let us introduce a system with a set of operable states $j = 0 \dots n$ and one failed state F. State 0 is an initial fully operable state and states $j = 1 \dots n$ are degraded operable states. In degraded states, the system is operable, however for larger j , the system is closer to the failed state, and system functioning performance may be lower.

Also, let us introduce the following scheme of transitions between the system states:

- From operable state $j = 0 \dots n$, the system can pass to the next more degraded operable state, $j + 1$, if $j < n$, or to failed state F, if $j = n$, with given failure rate λ_j .
- From the operable state, $j = 1 \dots n$, the system can pass back to the less degraded operable state, $j - 1$, with given repair rate μ_j .
- In each operable state, $j = 0 \dots n$, some kind of critical failure may occur in the system, which transfers it from the operable state directly to failed state F with the given rate of σ_j .
- In the failed state, a full restore procedure for the system is used, which returns the system to initial operable state 0 with given rate γ .

Now, let us introduce Markov chain (figure 1), which represents graphically the reliability model, described above:

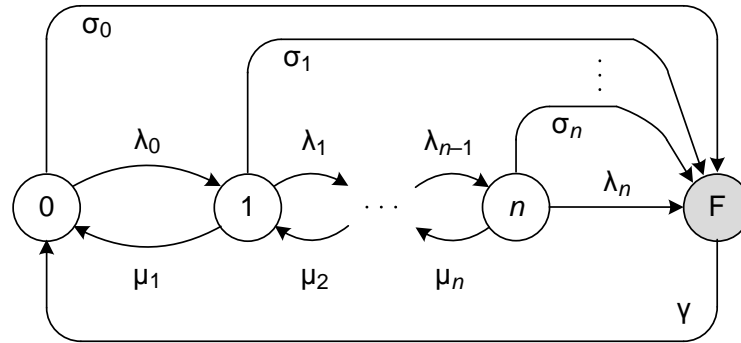


Figure 1. Specialized Markov chain for the reliability model of the system with dual-type failures, non-synchronized repairs and full restore after reaching the failed-state.

Accordingly, the system of Kolmogorov-Chapman equations for the stationary case is as follows:

$$\left\{ \begin{array}{l} P_0 + P_1 + \dots + P_n + P_F = 1; \\ -(\lambda_0 + \sigma_0)P_0 + \mu_1 P_1 + \gamma P_F = 0; \\ \lambda_0 P_0 - (\mu_1 + \lambda_1 + \sigma_1)P_1 + \mu_2 P_2 = 0; \\ \vdots \\ \lambda_{n-2} P_{n-2} - (\mu_{n-1} + \lambda_{n-1} + \sigma_{n-1})P_{n-1} + \mu_n P_n = 0; \\ \lambda_{n-1} P_{n-1} - (\mu_n + \lambda_n + \sigma_n)P_n = 0; \\ \sigma_0 P_0 + \dots + \sigma_{n-1} P_{n-1} + (\lambda_n + \sigma_n)P_n - \gamma P_F = 0. \end{array} \right. \quad (1)$$

An availability factor of the system is equal to the sum of probabilities of all operable states:

$$K_S = \sum_{j=0}^n P_j.$$

Next, we can easily calculate the mean time to repair by taking into consideration the fact that the system from the failed state can transit only to the initial operable state with rate γ :

$$T_R = 1/\gamma.$$

Finally, mean time to failure can be derived by using the well-known in reliability theory identity for repairable systems $K_S = T_F / (T_F + T_R)$:

$$T_F = K_S / (\gamma(1 - K_S)).$$

So, to obtain the stationary availability factor and mean time to failure, we need to solve the equations system and get all stationary probabilities.

It should be noted, that solving equation system (1) is a multi-level iterative process with cubic computational complexity $\sim 2(n+2)^3$. However, the author obtained a general analytic solution (for any given $n > 0$) in the form of the following matrix formula:

$$\begin{aligned} \Psi &= \prod_{j=1}^n \begin{bmatrix} \lambda_{n+1-j} & 0 & 0 \\ 1 & \mu_{n+1-j} & \sigma_{n+1-j} \\ 1 & \mu_{n+1-j} & \lambda_{n+1-j} + \sigma_{n+1-j} \end{bmatrix}; \\ \begin{bmatrix} U & 0 & 0 \\ V & 0 & W \\ M & 0 & D \end{bmatrix} &= \Psi \times \begin{bmatrix} \lambda_0 & 0 & 0 \\ 1 & 0 & \sigma_0 \\ 1 & 0 & \lambda_0 + \sigma_0 \end{bmatrix}; \\ K_S &= \frac{\gamma M}{\gamma M + D}; \quad T_F = \frac{M}{D}; \quad T_R = \frac{1}{\gamma}. \end{aligned} \quad (2)$$

A matrix formula requires multiplication of the $n + 1$ square matrixes with sizes 3 by 3, which contains all source reliability parameters of the system. The matrix formula has linear computational complexity $\sim 138(n+1)$. In general case, the matrix formula produces a result matrix with five non-zero coefficients U , V , W , M and D , and only two of them, M and D , are used for calculation of the stationary availability factor and mean time to failure.

By taking into the consideration the associative feature of matrix multiplication $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$, the author obtained the following recurrent computational scheme, which has also linear, but significantly less computational complexity $\sim 36n$ and provides fast calculation of the stationary availability factor and mean time to failure:

$$\begin{cases} U^{(1)} = \lambda_0; \quad V^{(1)} = 1; \quad M^{(1)} = 1; \\ W^{(1)} = \sigma_0; \quad D^{(1)} = \lambda_0 + \sigma_0; \\ r = 1 \dots n; \\ \left\{ \begin{aligned} U^{(r+1)} &= \lambda_r U^{(r)}; \\ V^{(r+1)} &= \sigma_r M^{(r)} + \mu_r V^{(r)} + U^{(r)}; \\ M^{(r+1)} &= \lambda_r M^{(r)} + V^{(r+1)}; \\ W^{(r+1)} &= \sigma_r D^{(r)} + \mu_r W^{(r)}; \\ D^{(r+1)} &= \lambda_r D^{(r)} + W^{(r+1)}; \\ M &= M^{(n+1)}; \quad D = D^{(n+1)}; \end{aligned} \right. \\ K_S = \frac{\gamma M}{\gamma M + D}; \quad T_F = \frac{M}{D}; \quad T_R = \frac{1}{\gamma}. \end{cases} \quad (3)$$

3. Specialized Markov chain in the reliability model of the dual-level disk array 'RAID-10'

The fault-tolerant dual-level disk array 'RAID-10' is a data storage system, which combines high reliability in the inner level due to the data mirroring in 'RAID-1' technology and high-performance in the outer level due to data striping in 'RAID-0' technology. In the inner level, there is a set of n independent dual-disk 'RAID-1' arrays. Data in each dual-disk 'RAID-1' array are always

synchronized between two disks. In the outer level, a set of n dual-disk 'RAID-1' arrays are united to an array of arrays by using the 'RAID-0' technology (figure 2). Data in the 'RAID-0' array are distributed between blocks, located on the different 'RAID-1' arrays.

The 'RAID-10' array in total contains $2n$ disks ($2n \geq 4$). User data capacity is equal to 50% of total disk space because of data mirroring between two disks in each 'RAID-1' array.

In the best case, 'RAID-10' is operable even if n disks fail and all of them belong to different 'RAID-1' arrays, and there is no pair of disks, which belong to the same 'RAID-1' array. In the worst case, even failure of two disks can cause the whole 'RAID-10' failure and data loss, if both of them belong to the same 'RAID-1' array, because the 'RAID-0' technology does not provide any fault-tolerance and failure of any dual-disk 'RAID-1' array, which causes failure of the whole 'RAID-10'.

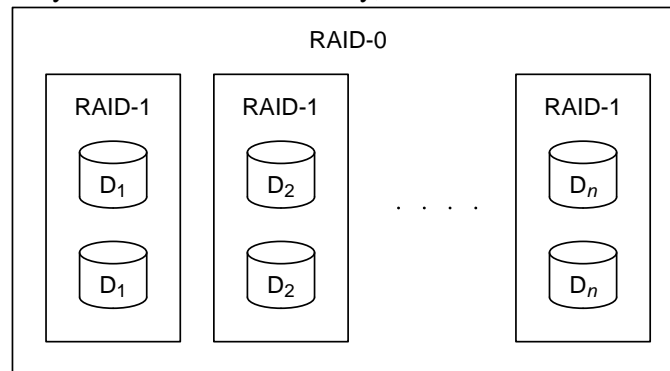


Figure 2. The structure of the dual-level disk array 'RAID-10'.

Let us now overview the offered reliability model of the data storage system, based on the dual-level disk array 'RAID-10' with $2n$ identical disks (the same capacity, manufacturer and model of all disks). The disk failure rate is λ . Disks can fail independently.

The initial state of the data storage system is 0 (all disks are operable). From this state, in case of failure of one of the $2n$ disks, the system passes to state 1 (one of non-fatal failures occurs).

From state 1, in case of failure of one of the $2n - 1$ disks, the system passes either to state 2 (two of non-fatal failures occurred) if the failed disk is the one of the $2n - 2$ disks, which do not belong to the 'RAID-1' array with the previously failed disk, or to failed state F if currently and previously failed disks belong to the same 'RAID-1' array.

From state 2, in case of failure of one of the $2n - 2$ disks, the system passes either to state 3 (three of non-fatal failures occurred) if the failed disk is the one of the $2n - 4$ disks, which do not belong to the one of 'RAID-1' arrays with one of the previously failed disks, or to failed state F, if current and one of the previously failed disks belong to the same 'RAID-1'.

Accordingly, in state n (n of non-fatal failures occurred), in each 'RAID-1' array, one of two disks will be failed. Thus, in state n , failure of any remaining n disks will transfer the system to failed state F because, in this case, one of the 'RAID-1' arrays will fail.

Next, let us assume, that 'RAID-1' arrays are independent for repairs, and in each of them, the rebuild process can be started independently after replacement of the failed disk. Completions of the rebuild processes in different 'RAID-1' arrays are also independent and not synchronized. The 'RAID-1' array rebuild rate is μ . We must mention that for simplification of the reliability model, we will consider the disk replacement time as negligible (assuming that the unlimited amount of the hot-spare disks is available). Completion of the rebuild process in the separate 'RAID-1' array returns the system from state $j = 1 \dots n$ back to state $j - 1$.

Also it should be noted that for the rebuild process in each of the 'RAID-1' arrays, we will take into consideration additional disk read error rate ε . It should be added to the failure rate of the disk, from which data are copied to the replaced disk in the 'RAID-1' array during the rebuild process. Failure of the data source disk as well as the read error on it during the rebuild process can cause failure of the whole 'RAID-1' array and transfer the system to failed state F.

Next, let us take into consideration critical array controller errors with rate σ , which transfer the system from any operable state $j = 0 \dots n$ directly to failed state F.

Finally, let us assume that the system does not include any additional hardware and software to provide total backup of the user data. In case of reaching the failed state, all data will be lost. So, we will consider the rate of the system full recovery rate, which is zero, $\gamma = 0$.

Now, we may introduce Markov chain (figure 3), which represents graphically the reliability model, described above:

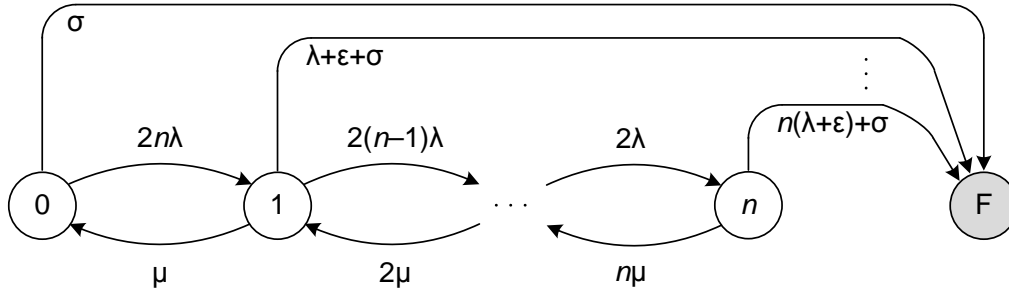


Figure. 3. The specialized reliability model for the dual-level disk array 'RAID-10'.

This model is a particular case of the specialized Markov chain, discussed above, with the following substitutions of the source reliability parameters:

$$\begin{cases} \lambda_0 = 2n\lambda; & \lambda_j = 2(n-j)\lambda; \\ \sigma_0 = \sigma; & \sigma_j = \sigma + j(\lambda + \varepsilon); \\ \gamma = 0; & \mu_j = j\mu; \\ & j = 1 \dots n. \end{cases} \quad (4)$$

As the system full recovery rate is zero, $\gamma = 0$, the stationary availability factor is zero and mean time to recovery is infinite for the discussed 'RAID-10' array. However, mean time to data loss for 'RAID-10' is equal to the mean time to failure in specialized Markov chain. So, after substitution of the source reliability parameters into the above-mentioned recurrent calculation scheme (3), we obtain the following recurrent scheme for calculation of the mean time to data loss of the 'RAID-10' array:

$$\begin{cases} U^{(1)} = 2n\lambda; & V^{(1)} = 1; & M^{(1)} = 1; \\ W^{(1)} = \sigma; & D^{(1)} = 2n\lambda + \sigma; \\ & r = 1 \dots n; \\ \left\{ \begin{array}{l} U^{(r+1)} = 2(n-r)\lambda U^{(r)}; \\ V^{(r+1)} = (\sigma + r(\lambda + \varepsilon))M^{(r)} + r\mu V^{(r)} + U^{(r)}; \\ M^{(r+1)} = 2(n-r)\lambda M^{(r)} + V^{(r+1)}; \\ W^{(r+1)} = (\sigma + r(\lambda + \varepsilon))D^{(r)} + r\mu W^{(r)}; \\ D^{(r+1)} = 2(n-r)\lambda D^{(r)} + W^{(r+1)}; \\ M = M^{(n+1)}; & D = D^{(n+1)}; \\ T_{DL} = M / D. \end{array} \right. \end{cases} \quad (5)$$

4. The calculation example for the mean time to data loss of the 'RAID-10' array

A 'RAID-10' array with total number of disks $2n$ is given. The disk failure rate is $\lambda = 1/120000 \text{ hour}^{-1}$, the additional disk read error rate during the rebuild process is $\varepsilon = 1/112 \text{ hour}^{-1}$, the 'RAID-1' array rebuild rate is $\mu = 1/9 \text{ hour}^{-1}$ and the array controller critical error rate is $\sigma = 1/1200000 \text{ hour}^{-1}$. Let us calculate the mean time to data loss of the 'RAID-10' array for the total number of disks $2n = 2, 4, 6, 8, 10$ and 12 .

Calculations by the recurrent scheme (5) give us the following values, presented in table 1:

Table 1. Mean times to data loss of the ‘RAID-10’ array for different total numbers of disks.

$2n$	T_{DL} (hours)
4	302642
6	220282
8	173159
10	142645
12	121275
14	105473
16	93315

Calculation results show that the mean time to data loss of the ‘RAID-10’ array is moderately decreased with the increase in the total number of disks.

5. Conclusion

In this scientific paper, a specialized Markov chain for the reliability model of the fault-tolerant system with dual-type failures, non-synchronized repairs and full restore after reaching the failed state is discussed. A generalized iterative procedure for calculation of the system’s stationary availability factor, mean time to repair and mean time of failure is also introduced. Finally, an application of the specialized Markov chain in the reliability model of the fault-tolerant dual-level disk arrays ‘RAID-10’ and calculation of the mean time to data loss by using the generalized iterative procedure is also observed.

Scientific results were used by the author in designing of fault-tolerant data storage systems based on the ‘RAID-10’ arrays for Moscow Power Engineering Institute, Nuclear Power Plant “Balakovo” and several other enterprises.

Acknowledgements

The author thanks professor I I Ladygin, Moscow Power Engineering Institute (MPEI), for scientific support, and S N Khorkov, chief network administrator of MPEI, for technical support.

References

- [1] Cherkesov G N 2005 *Reliability of Hardware and Software Systems* (Saint-Petersburg: Piter)
- [2] Polovko A M and Gurov S V 2006 *Basis of Reliability Theory* (Saint-Petersburg: BHV-Petersburg)
- [3] Shooman M L 2002 *Reliability of computer systems and networks* (John Wiley & Sons, Inc.)
- [4] Koren I and Krishna C M 2007 *Fault-Tolerant Systems* (Morgan Kaufmann Publishers)
- [5] Anderson W J 1991 *Continuous-Time Markov Chains* (Springer-Verlag)
- [6] Wang Z, Wang T and Yang X 1992 *Birth and death processes and Markov chains* (Springer-Verlag)
- [7] Schmidt K 2006 *High-availability and disaster recovery* (Springer-Verlag)
- [8] Nelson S 2011 *Pro data backup and recovery* (Apress)
- [9] Rahman P A, Muraveva E A and Sharipov M I 2016 *Key Eng. Mat.* **685** 805-810
- [10] Rahman P A and Bobkova E Yu 2016 *IOP Conf. Ser.: Mater. Sci. Eng.* (Tomsk) vol 124 (IOP Publishing) 012023