

Mean Value Analysis

Raj Jain

Washington University in Saint Louis
Jain@eecs.berkeley.edu or Jain@wustl.edu

A Mini-Course offered at UC Berkeley, Sept-Oct 2012

These slides and audio/video recordings are available on-line at:

<http://amplab.cs.berkeley.edu/courses/queue>

and <http://www.cse.wustl.edu/~jain/queue>

UC Berkeley, Fall 2012

©2012 Raj Jain

34-1



- Exact solution using an iterative method with several assumptions
- Key steps
- Assumption

UC Berkeley, Fall 2012

©2012 Raj Jain

34-2

Mean-Value Analysis (MVA)

- Mean-value analysis (MVA) allows solving closed queueing networks
- It gives the mean performance.
The variance computation is not possible using this technique.
- Initially limit to fixed-capacity service centers and delay centers.

4 Steps:

1. Given a closed queueing network with N jobs:

$$R_i(N) = S_i (1 + Q_i(N-1))$$

- Here, $Q_i(N-1)$ is the mean queue length at i^{th} device with $N-1$ jobs in the network.

- It assumes that the service is memoryless.

Note: This is not PASTA. Arrivals are not Poisson.



UC Berkeley, Fall 2012

©2012 Raj Jain

34-3

Mean-Value Analysis (MVA)

- Since the performance with no users ($N=0$) can be easily computed, performance for any number of users can be computed iteratively.
2. Given the response times at individual devices, the system response time using the general response time law is:

$$R(N) = \sum_{i=1}^M V_i R_i(N)$$

3. The system throughput using the interactive response time law is:

$$X(N) = \frac{N}{R(N) + Z}$$

UC Berkeley, Fall 2012

©2012 Raj Jain

34-4

Mean-Value Analysis (MVA)

- The device throughputs measured in terms of jobs per second are:

$$X_i(N) = X(N) V_i$$

- The device queue lengths with N jobs in the network using Little's law are:

$$Q_i(N) = X_i(N) R_i(N) = X(N) V_i R_i(N)$$

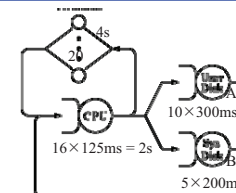
- Response time equation for delay centers is simply:

$$R_i(1) = S_i$$

- Earlier equations for device throughputs and queue lengths apply to delay centers as well.

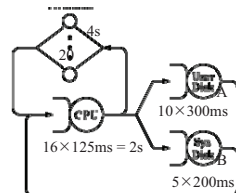
$$Q_i(0) = 0$$

Example 34.2



- Consider a timesharing system
- Each user request makes ten I/O requests to disk A, and five I/O requests to disk B.
- The service times per visit to disk A and disk B are 300 and 200 milliseconds, respectively.
- Each request takes two seconds of CPU time and the user think time is four seconds.
 $S_A = 0.3, V_A = 10 \Rightarrow D_A = 3$
 $S_B = 0.2, V_B = 5 \Rightarrow D_B = 1$
 $D_{CPU} = 2, V_{CPU} = V_A + V_B + 1 = 16 \Rightarrow S_{CPU} = 0.125$
 $Z = 4, \text{ and } N = 20$

Example 34.2 (Cont)



- Initialization:

- Number of users: $N=0$
- Device queue lengths: $Q_{CPU}=0, Q_A=0, Q_B=0$

- Iteration 1: Number of users: $N=1$

- Device response times:

$$R_{CPU} = S_{CPU}(1 + Q_{CPU}) = 0.125(1 + 0) = 0.125$$

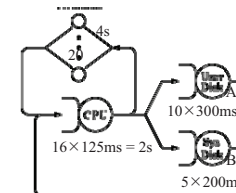
$$R_A = S_A(1 + Q_A) = 0.3(1 + 0) = 0.3$$

$$R_B = S_B(1 + Q_B) = 0.2(1 + 0) = 0.2$$

- System Response time:

$$\begin{aligned} R &= R_{CPU}V_{CPU} + R_A V_A + R_B V_B \\ &= 0.125 \times 16 + 0.3 \times 10 + 0.2 \times 5 = 6 \end{aligned}$$

Example 34.2 (Cont)



- System Throughput:

$$X = N/(R + Z) = 1/(6 + 4) = 0.1$$

- Device queue lengths:

$$Q_{CPU} = X R_{CPU} V_{CPU} = 0.1 \times 0.125 \times 16 = 0.2$$

$$Q_A = X R_A V_A = 0.1 \times 0.3 \times 10 = 0.3$$

$$Q_B = X R_B V_B = 0.1 \times 0.2 \times 5 = 0.1$$

- Iteration 2: Number of users: $N=2$

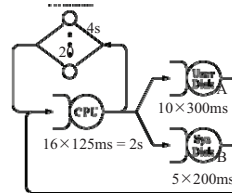
- Device response times:

$$R_{CPU} = S_{CPU}(1 + Q_{CPU}) = 0.125(1 + 0.2) = 0.15$$

$$R_A = S_A(1 + Q_A) = 0.3(1 + 0.3) = 0.39$$

$$R_B = S_B(1 + Q_B) = 0.2(1 + 0.1) = 0.22$$

Example 34.2 (Cont)



2. System Response time:

$$R = R_{CPU}V_{CPU} + R_A V_A + R_B V_B$$

$$= 0.15 \times 16 + 0.39 \times 10 + 0.22 \times 5 = 7.4$$

3. System Throughput: $X = N/(R+Z) = 2/(7.4+4) = 0.175$

4. Device queue lengths:

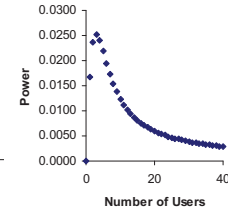
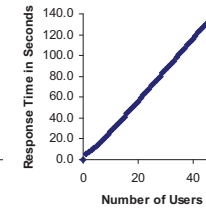
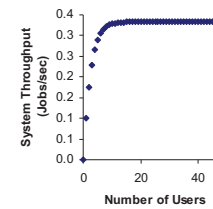
$$Q_{CPU} = X R_{CPU} V_{CPU} = 0.175 \times 0.15 \times 16 = 0.421$$

$$Q_A = X R_A V_A = 0.175 \times 0.39 \times 10 = 0.684$$

$$Q_B = X R_B V_B = 0.175 \times 0.22 \times 5 = 0.193$$

MVA Results for Example 34.2

Iteration #	Response Time				System Throughput	Queue Lengths		
	CPU	Disk A	Disk B	System		CPU	Disk A	Disk B
1	0.125	0.300	0.200	6.000	0.100	0.200	0.300	0.100
2	0.150	0.390	0.220	7.400	0.175	0.421	0.684	0.193
3	0.178	0.505	0.239	9.088	0.229	0.651	1.158	0.273
4	0.206	0.647	0.255	11.051	0.266	0.878	1.721	0.338
5	0.235	0.816	0.268	13.256	0.290	1.088	2.365	0.388
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
17	0.370	3.962	0.300	47.045	0.333	1.974	13.195	0.499
18	0.372	4.259	0.300	50.032	0.333	1.981	14.187	0.499
19	0.373	4.556	0.300	53.022	0.333	1.987	15.181	0.500
20	0.373	4.854	0.300	56.016	0.333	1.991	16.177	0.500



Box 34.1: MVA Algorithms

Inputs:	Outputs:
N = number of users	X = system throughput
Z = think time	Q_i = average # of jobs at i th device
M = number of devices	R_i = response time of i th device
S_i = service time/visit to i th device	R = system response time
V_i = number of visits to i th device	U_i = utilization of the i th device

Initialization: FOR $i = 1$ TO M DO $Q_i = 0$

Iterations:

FOR $n = 1$ TO N DO
BEGIN

FOR $i = 1$ TO M DO $R_i = \begin{cases} S_i(1 + Q_i) & \text{Fixed capacity} \\ S_i & \text{Delay centers} \end{cases}$

$$R = \sum_{i=1}^M R_i V_i$$

$$X = \frac{N}{Z+R}$$

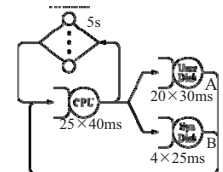
FOR $i = 1$ TO M DO $Q_i = X V_i R_i$

END

Device throughputs: $X_i = X V_i$

Device utilizations: $U_i = X S_i V_i$

Quiz 34A: MVA



Part 1: Fill in the rows for $N=0$ and $N=1$ only.

$$R_i = S_i(1 + Q_i)$$

$$R = \sum_{i=1}^M R_i V_i$$

$$X = \frac{N}{Z+R}$$

$$Q_i = X V_i R_i$$

V_i	25	20	4					$Z=5$				
S_i	0.04	0.03	0.025									
N	R_C	R_A	R_B	$V_C R_C$	$V_A R_A$	$V_B R_B$	R	$R+Z$	X	Q_C	Q_A	Q_B
0	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____
1	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____
2	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____	_____

Part 2: Fill in the row for $N=2$.

MVA Assumptions

- MVA is applicable only if the network is a product form network with exponentially distributed service times.
 1. **Job flow balance:** # In = # out \Rightarrow No buffer overflow
 2. **One step behavior:** Only one job in or out at a time \Rightarrow No bulk arrivals or service
 3. Only **fixed-capacity service centers or delay centers**
Load dependent servers *can be* included but not covered here.
 4. **Exponentially distributed service times** for all centers
 5. **Device Homogeneity:** A device's service rate for a particular class does not depend on the state of the system in any way except for the total device queue length and the designated class's queue length.

UC Berkeley, Fall 2012

©2012 Raj Jain

34-15

MVA Assumptions (Cont)

Device homogeneity implies the following:

- a. **Single Resource Possession:** A job may not be present (waiting for service or receiving service) at two or more devices at the same time.
- b. **No Blocking:** A device renders service whenever jobs are present; its ability to render service is not controlled by any other device.
- c. **Independent Job Behavior:** Interaction among jobs is limited to queueing for physical devices, for example, there should not be any synchronization requirements.
- d. **Local Information:** A device's service rate depends only on local queue length and not on the state of the rest of the system.



UC Berkeley, Fall 2012

©2012 Raj Jain

34-16

MVA Assumptions (Cont)

- e. **Fair Service:** If service rates differ by class, the service rate for a class depends only on the queue length of that class at the device and not on the queue lengths of other classes. This means that the servers do not discriminate against jobs in a class depending on the queue lengths of other classes. (No priority)

UC Berkeley, Fall 2012

©2012 Raj Jain

34-17

Summary



1. MVA allows exact analysis of closed queueing networks.
Given performance of N-1 users, get performance for N users.
2. 4 Steps:
$$R_i = S_i(1 + Q_i)$$
$$R = \sum_{i=1}^M R_i V_i$$
$$X = \frac{N}{Z + R}$$
$$Q_i = X V_i R_i$$
3. Assumptions: Exponential service times, flow balance, one-step behavior, device homogeneity

UC Berkeley, Fall 2012

©2012 Raj Jain

34-18

Review of Key Concepts

1. Kendall Notation: A/S/m/B/k/SD, M/M/1
2. Little's Law:
Mean number in system = Arrival rate \times Mean time in system
3. Processes: Markov \Rightarrow Only one state required,
Poisson \Rightarrow IID and exponential inter-arrival
4. Operational Laws: No loss

Utilization Law:	$U_i = X_i S_i = X D_i$
Forced Flow Law:	$X_i = X V_i$
Little's Law:	$Q_i = X_i R_i$
General Response Time Law:	$R = \sum_{i=1}^M R_i V_i$
Interactive Response Time Law:	$R = \frac{N}{X} - Z$
Asymptotic Bounds:	$R \geq \max\{D, N D_{max} - Z\}$ $X \leq \min\{1/D_{max}, N/(D + Z)\}$

5. Mean Value Analysis: Single arrivals/service, no loss, exponential service time, device homogeneity

$$R_i = S_i(1 + Q_i) \quad R = \sum_{i=1}^M R_i V_i \quad X = \frac{N}{Z+R} \quad Q_i = X V_i R_i$$

Quiz 1: Post Quiz

True or False?

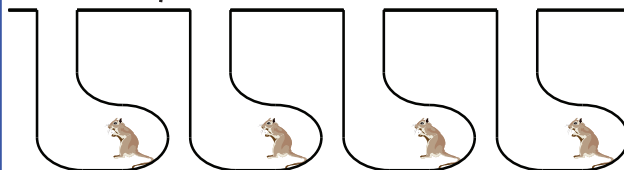
T F

- M/M/1/3/100 queue has 3 servers
- A single server queue with arrival rate of 1 jobs/sec and a service time of 0.5 seconds has server utilization of 0.5
- The delay in an G/G/ ∞ system is equal to the job service time.
- In a product form queueing network, the probability of a state can be obtained by multiplying state probabilities of individual queues.
- During a 10 second observation period, 400 jobs were serviced by a processor which can process 200 jobs per second. The processor utilization is 50%.
- MVA can be used to compute response times for non-product form networks.

Marks = Correct Answers _____ - Incorrect Answers _____ = _____

<http://amplab.cs.berkeley.edu/courses/queue/quiz1.html>

Performance Analysis Rat Holes



Workload Metrics Configuration Details

- Workload: Does not exercise the bottleneck, component under study, or the parameter.
- Metrics: Incomplete or wrong level
- Configuration: No experimental design
- Details: No validation

Reasons for not Accepting an Analysis

- This needs more analysis.
- You need a better understanding of the workload.
- It improves performance only for long IOs/packets/jobs/files, and most of the IOs/packets/jobs/files are short.
- It improves performance only for short IOs/packets/jobs/files, but who cares for the performance of short IOs/packets/jobs/files, its the long ones that impact the system.
- It needs too much memory/CPU/bandwidth and memory/CPU/bandwidth isn't free.
- It only saves us memory/CPU/bandwidth and memory/CPU/bandwidth is cheap.

See Box 10.2 on page 162 of the book for a complete list

Three Rules of Validation

- ❑ Do not trust the results of a **simulation model** until they have been validated by analytical modeling or measurements.
- ❑ Do not trust the results of an **analytical model** until they have been validated by a simulation model or measurements.
- ❑ Do not trust the results of a **measurement** until they have been validated by simulation or analytical modeling.

Experimental Design: Latex vs. troff

Factors and Levels

	Factor	-Level	+Level
A	Program	Latex	troff-me
B	Bytes	2100	25000
C	Equations	0	10
D	Floats	0	10
E	Tables	0	10
F	Footnotes	0	10

- ❑ 5 factors each at 2 levels $\Rightarrow 2^5$ experiments
- ❑ $2^{5-2} = 8$ experiments \Rightarrow which parameters are more important
- ❑ Run 2nd phase with smaller number of parameters and more levels.