# Decision Support System for Healthcare Capacity Planning

Canan Gunes Corlu, PhD, John Maleyeff, PhD, Chenshu Yang, MS, Tianhuai Ma, MS, Yanting Shen, MS Candidate
Administrative Sciences, Metropolitan College, Boston University

## Introduction

Capacity management of hospital staff and other resources is an important challenge faced by a healthcare administrator. Because of the variation in service times and the inability to inventory services, capacity buffers are required to ensure reasonable waiting times for patients. Figure 1 shows a generic example that is applicable to any queuing system. As the server utilization increases, waiting time will increase in a pattern commonly known as a hockey stick.
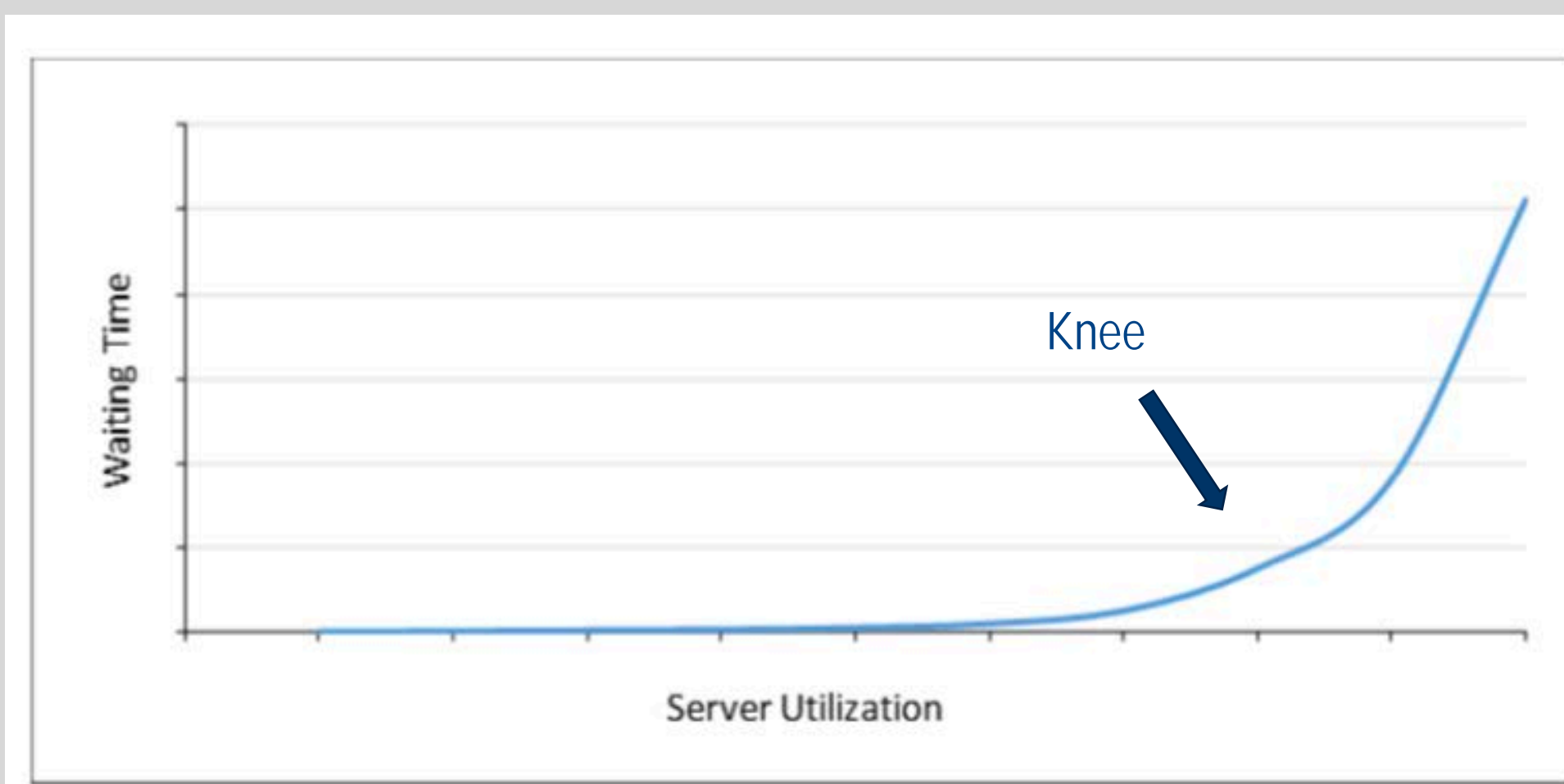


Figure 1  Server Utilization versus Waiting Time "Hockey-Stick."

A capacity plan will effectively balance the needs of the planner (i.e., by maximizing server utilization) and the needs of the patient (by minimizing waiting times). In some fields, the optimal system configuration takes place at a threshold referred to as the knee of the curve.

This study concerns the development of a decision support system (DSS) using Python to determine optimal capacity buffers using a Monte Carlo simulation (MCS) and a knee optimization model that allows for flexibility in specifying uncertain arrival patterns and service times.

## Assumptions

The queuing system is robust with the following structure:
- Single queue served by multiple parallel servers
- Infinite customer population
- Arrivals are deterministic or random
- Infinite queue size
- First-come first-served discipline
- Gamma distributed service time variation

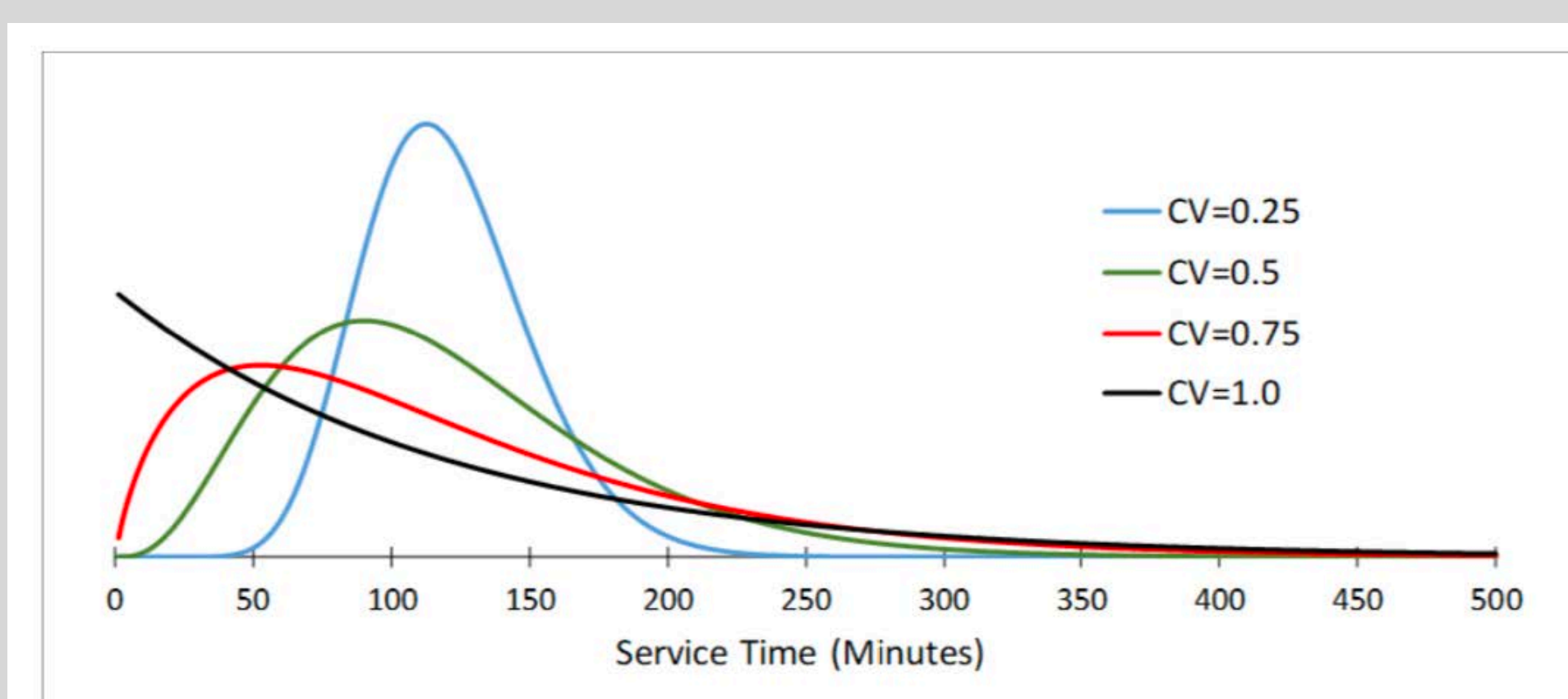Figure 2 shows gamma distributions for various values of the coefficient variation (CV).



Figure 2 Gamma Distributions (Average Service Time is 120 Minutes).

## Methods

**Simulation Model:**
- The performance statistics were generated based on the server utilization (ρ), which is the ratio of the customer arrival rate (λ) to the system's service rate.
- The inputs are the number of servers, CV, and the arrival pattern. The outputs are the customer's wait time in the system (Ws), wait time in the queue (Wq), the number of customers in the system (Ls) and the number of customers in the queue (Lq).

**Optimization Model:**
- This study uses a Kleinrock's power function to identify the level of ρ using the equation:

$$Power(\rho) = \frac{\rho}{\lambda\ Ws}$$

- The MCS finds the knee (optimal server utilization) by changing ρ from 40% to 95% (in increments of 5%) and simulating the system repeatedly over this range.

## Observations

**Congestion and Instability:**
- When the system is congested the variation of wait times increases as a percentage of the average wait time, and wait times behaved erratically.
- The variation of the time spent in a congested system exhibited a great deal of instability due to the significant autocorrelation among patient wait times.
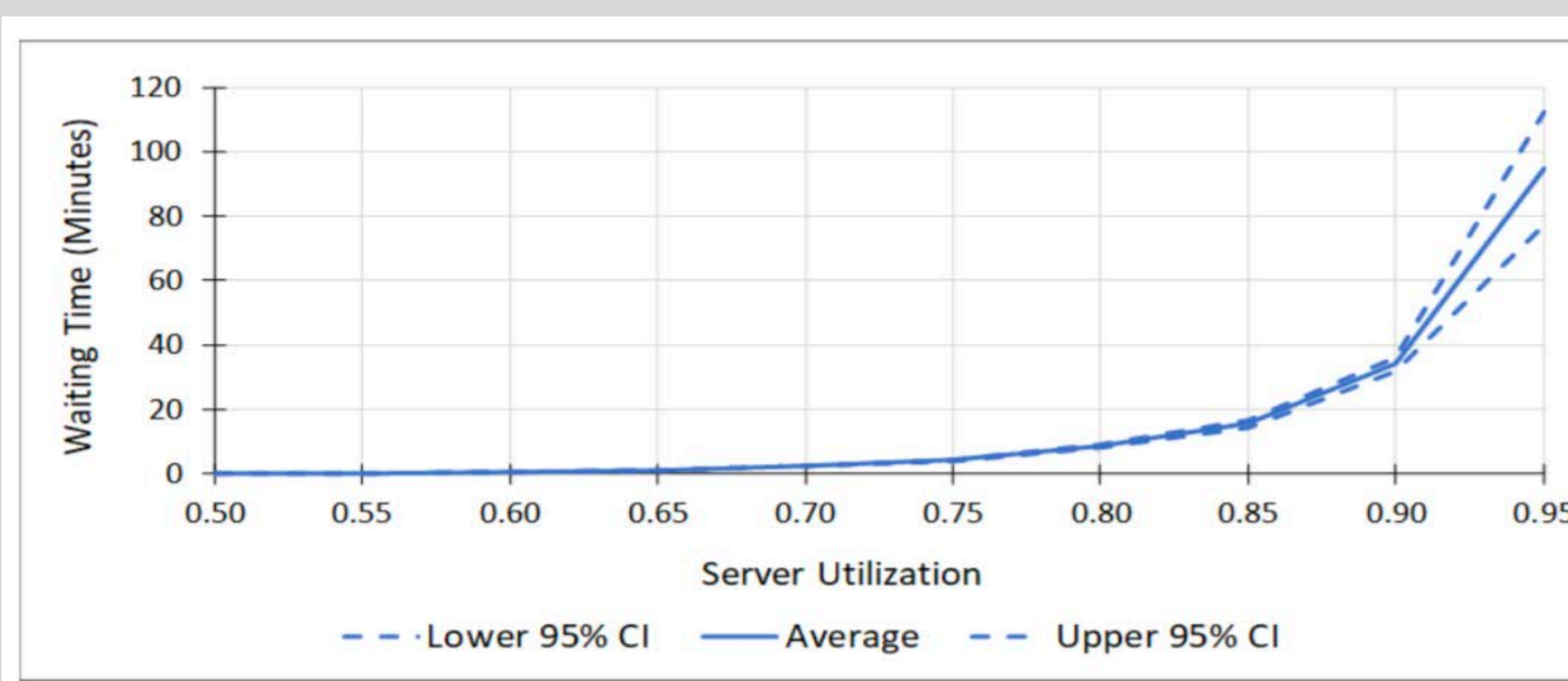- Increasing servers or reducing service times are the approaches to relieve congestion.



Figure 3 Confidence Intervals (s=15, CV=0.5, 6/hour random arrivals)

**Optimal Capacity Buffering:**
- A factorial experimental design was used to explore how the knee changed based on levels of key variables. Figure 4 shows the knee for 88 conditions.

| | Queuing System Assumptions | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Random Arrivals | | | | Scheduled Arrivals | | | |
| CV | 25% | 50% | 75% | 100% | 25% | 50% | 75% | 100% |
| 1 | 0.60 | 0.55 | 0.50 | 0.50 | 0.85 | 0.75 | 0.65 | 0.55 |
| 2 | 0.65 | 0.65 | 0.60 | 0.60 | 0.90 | 0.80 | 0.70 | 0.65 |
| 3 | 0.70 | 0.70 | 0.70 | 0.60 | 0.90 | 0.85 | 0.80 | 0.75 |
| 4 | 0.70 | 0.75 | 0.70 | 0.65 | 0.90 | 0.85 | 0.80 | 0.75 |
| 5 | 0.75 | 0.75 | 0.70 | 0.65 | 0.90 | 0.85 | 0.80 | 0.75 |
| 10 | 0.80 | 0.80 | 0.80 | 0.75 | 0.95 | 0.90 | 0.85 | 0.85 |
| 15 | 0.85 | 0.80 | 0.85 | 0.80 | 0.95 | 0.90 | 0.90 | 0.90 |
| 25 | 0.85 | 0.85 | 0.85 | 0.85 | 0.95 | 0.90 | 0.90 | 0.90 |
| 50 | 0.90 | 0.90 | 0.90 | 0.90 | 0.95 | 0.95 | 0.95 | 0.90 |
| 75 | 0.90 | 0.90 | 0.90 | 0.90 | 0.95 | 0.95 | 0.95 | 0.95 |
| 100 | 0.95 | 0.95 | 0.90 | 0.90 | 0.95 | 0.95 | 0.95 | 0.95 |

Figure 4 Optimal Knee Values

## Observations

- When determining the knee for each condition, 10 trials with 10000 iterations each were used to generate results for every level of server utilization.
- Figure 5 shows how the power function is used to determine the knee for 4 of the 88 conditions.
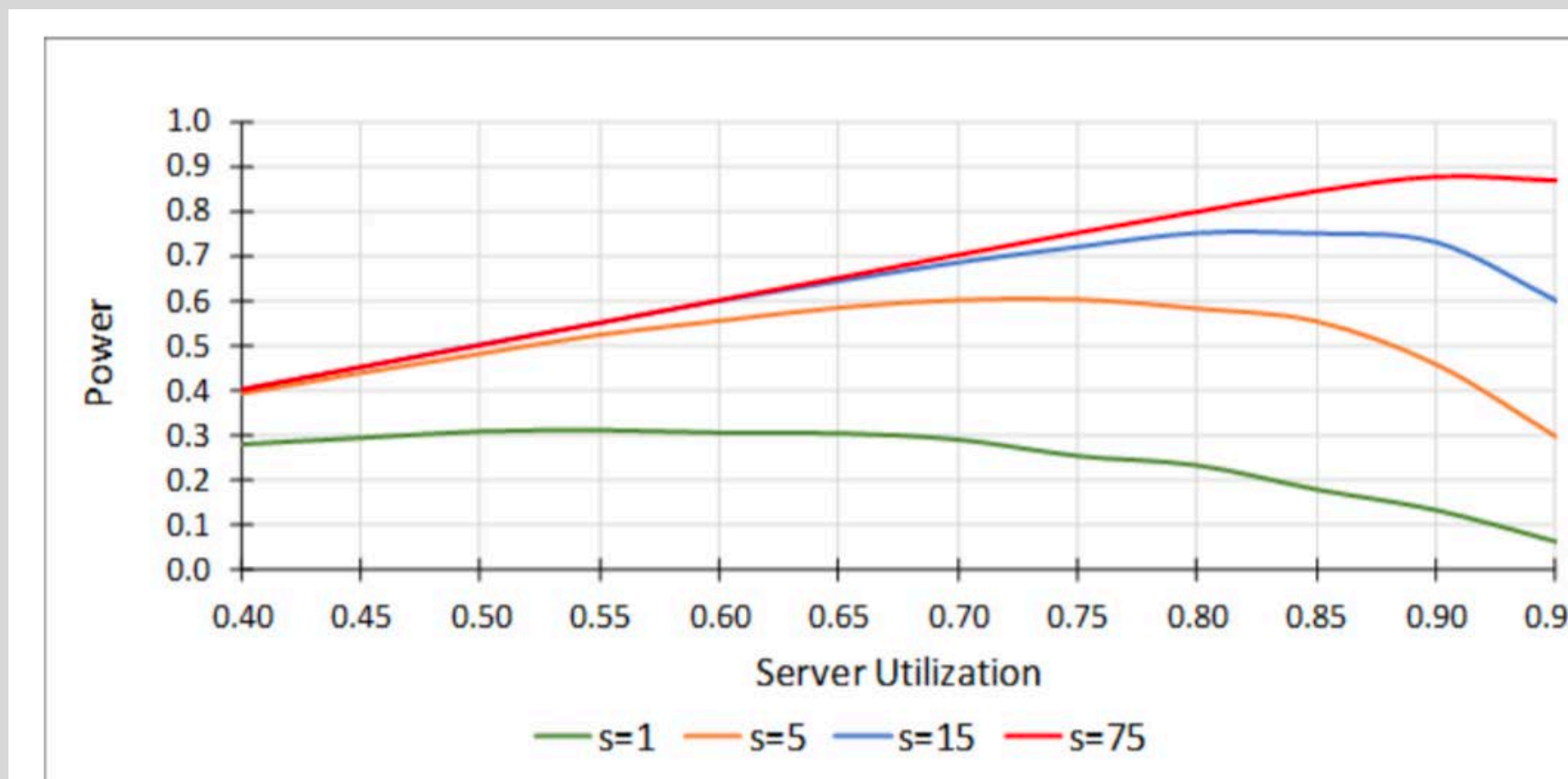


Figure 5 Utilization vs. Power Function (Random arrivals, CV=0.5).
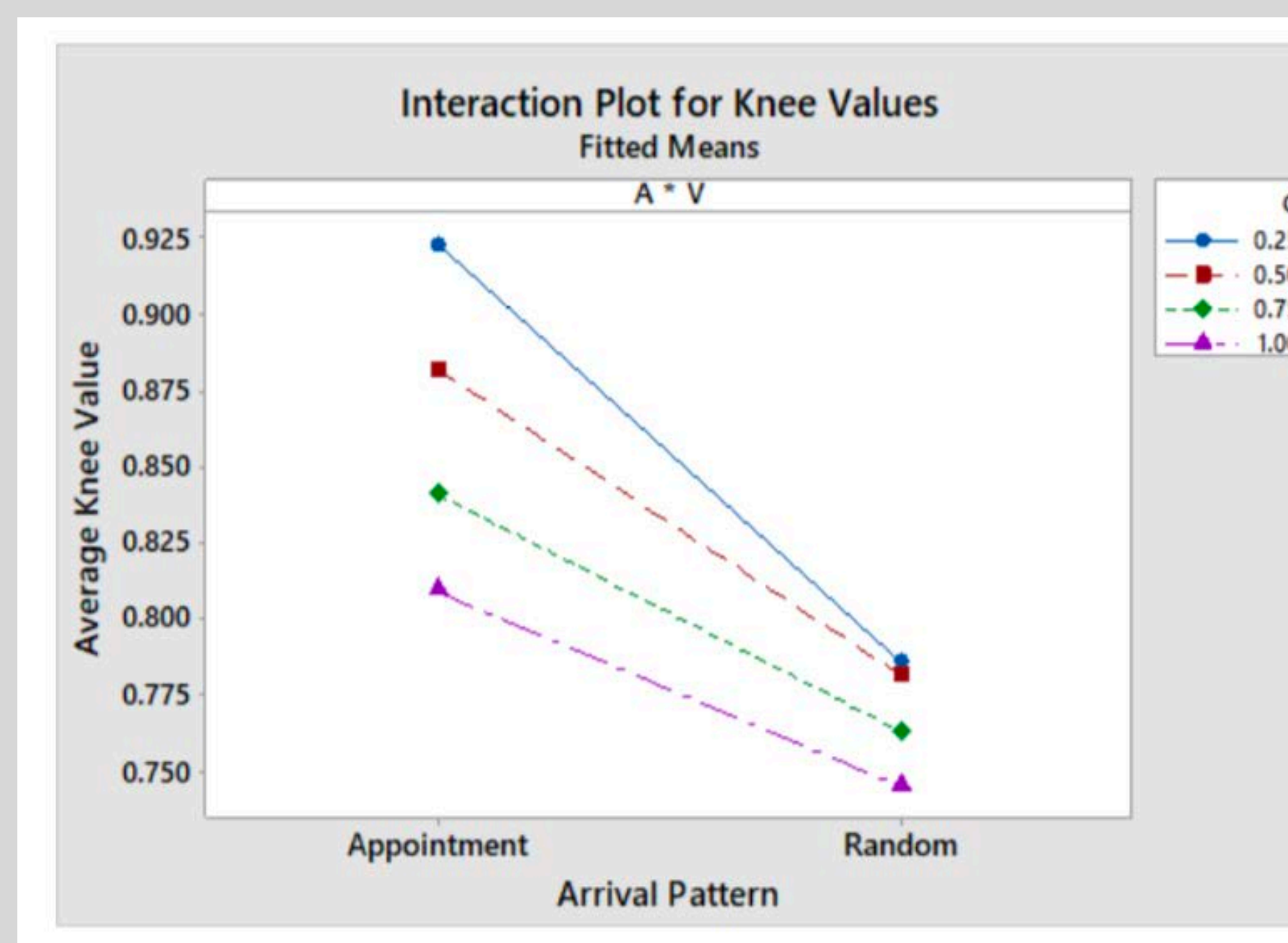- Figure 6 shows how the service time CV and the pattern of arrivals affected the knee.



Figure 6 Knee Values versus Arrival Patterns.

## Decision Support System - Input

A DSS has been developed in Python. Figure 7 shows the user interface, where the user enters the average service rate, average service time, CV, the number of servers, and a choice of random or scheduled arrivals.
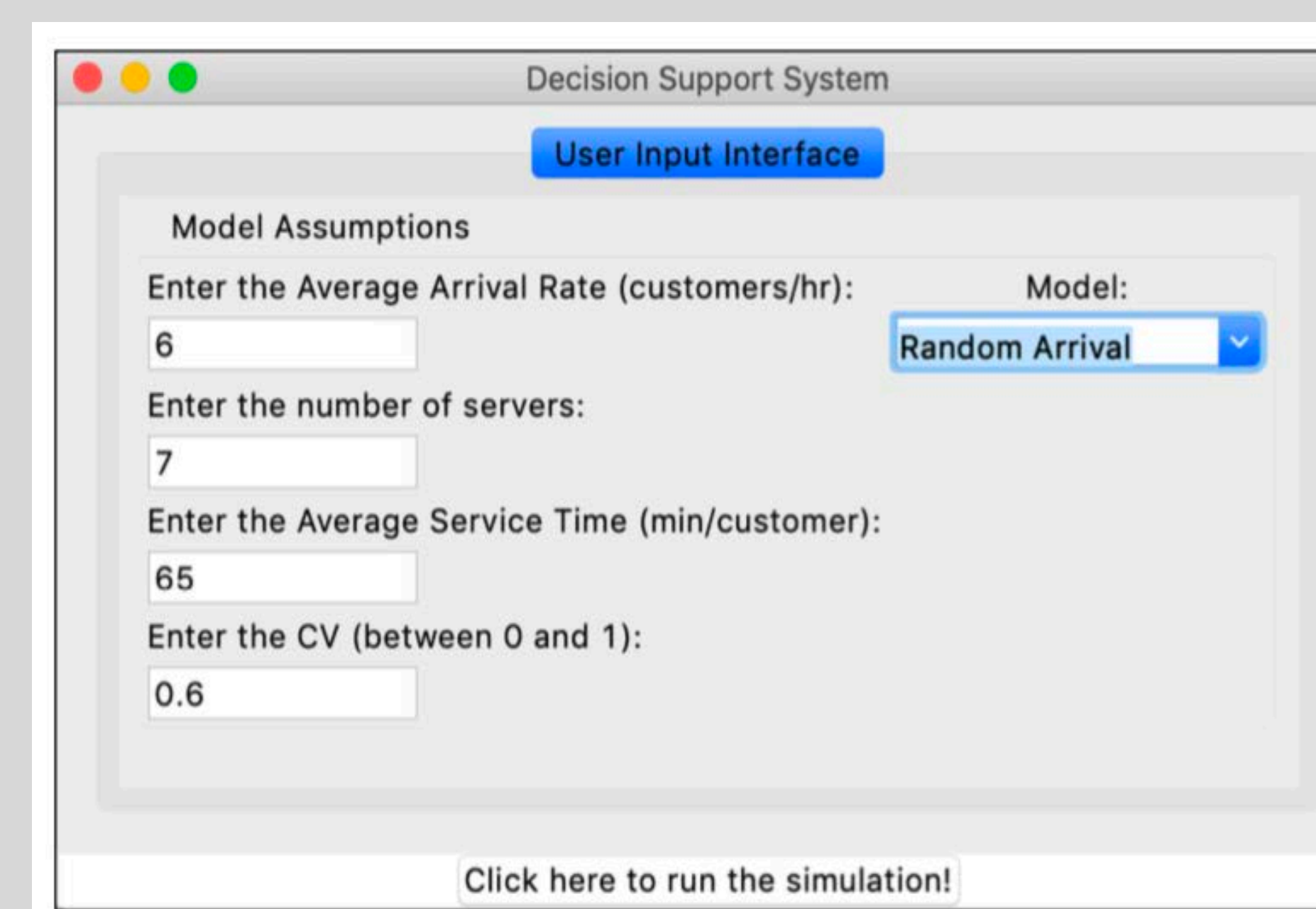


Figure 7 DSS User Interface

## Decision Support System - Output

The DSS begins by simulating the queuing system based on the user inputs, then it finds the knee using the MCS and knee optimization model.

Figure 8 shows the current server utilization compared to results with a range of utilization values based on service time changes, and an alternative analysis that focuses on changing the number of servers to achieve the optimal system configuration.
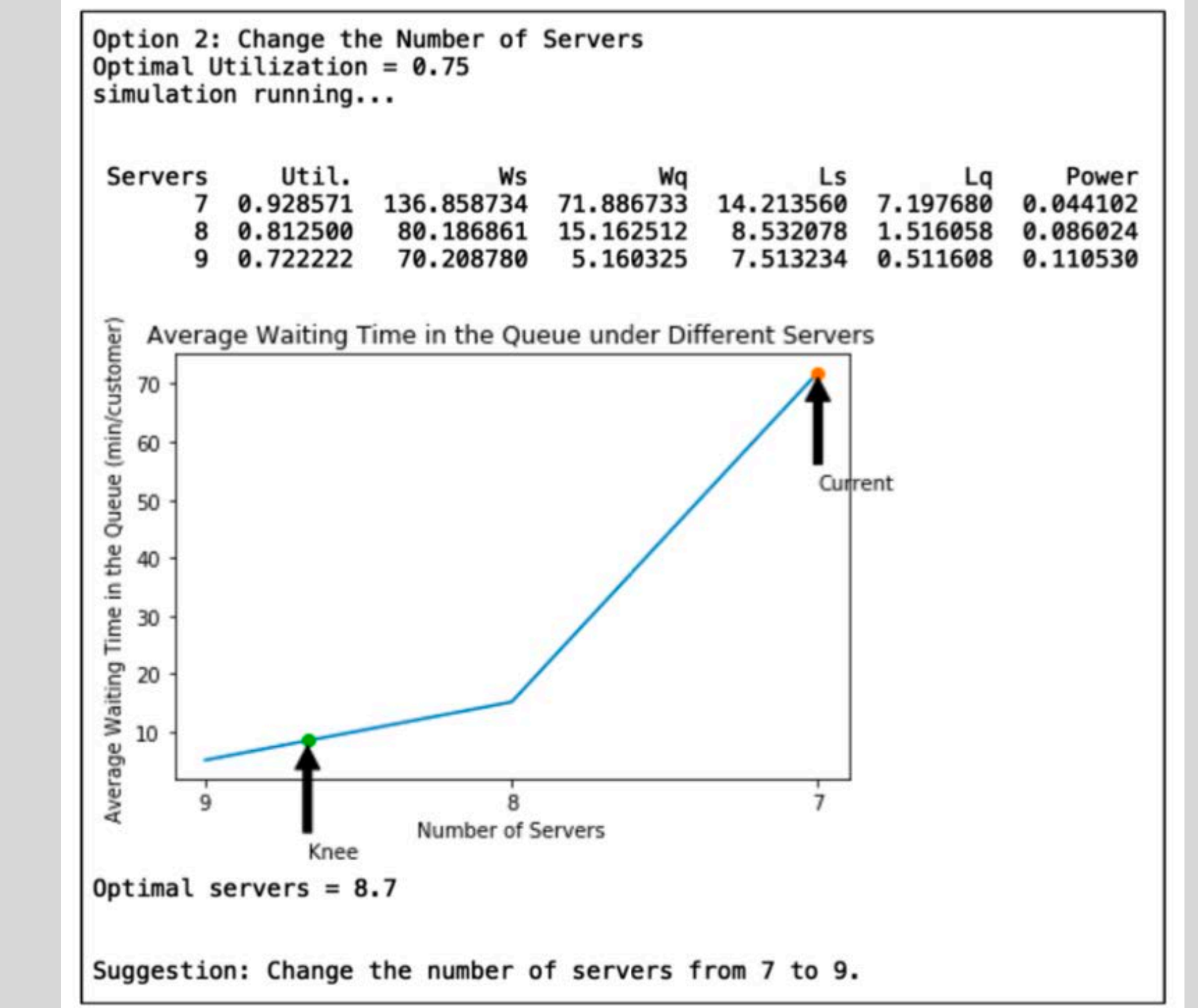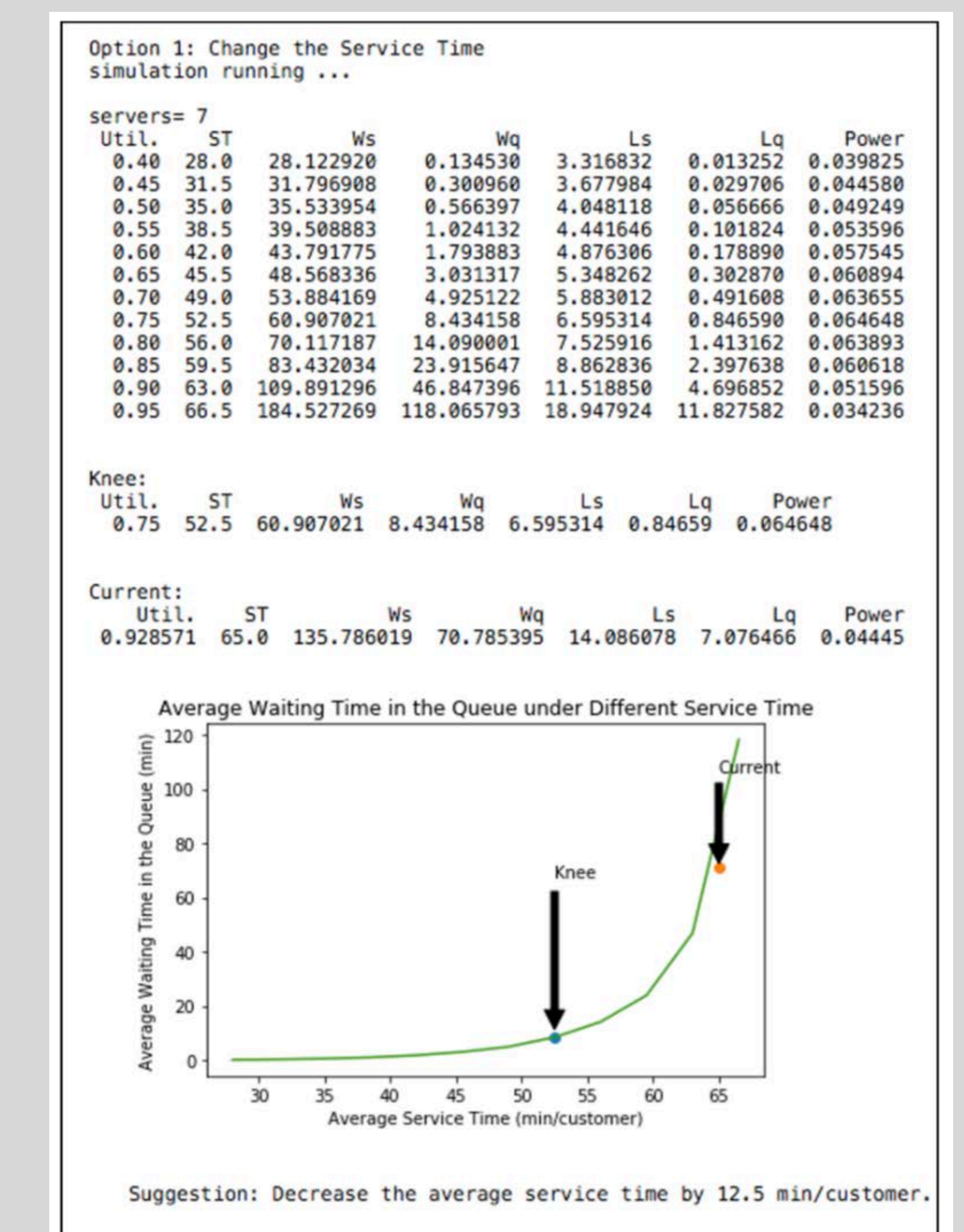


Figure 8 DSS Output Report