

Stochastic Optimization for Scheduling Service Systems

Brian Denton

Department of Industrial and Operations Engineering
University of Michigan
Ann Arbor, MI

February 20, 2024

Summary

Service System Scheduling Examples:

- Example 1: Single server scheduling
- Example 2: Multi-server scheduling
- Example 3: Bi-criteria scheduling of multi-server, multi-stage service system

Key Take Aways

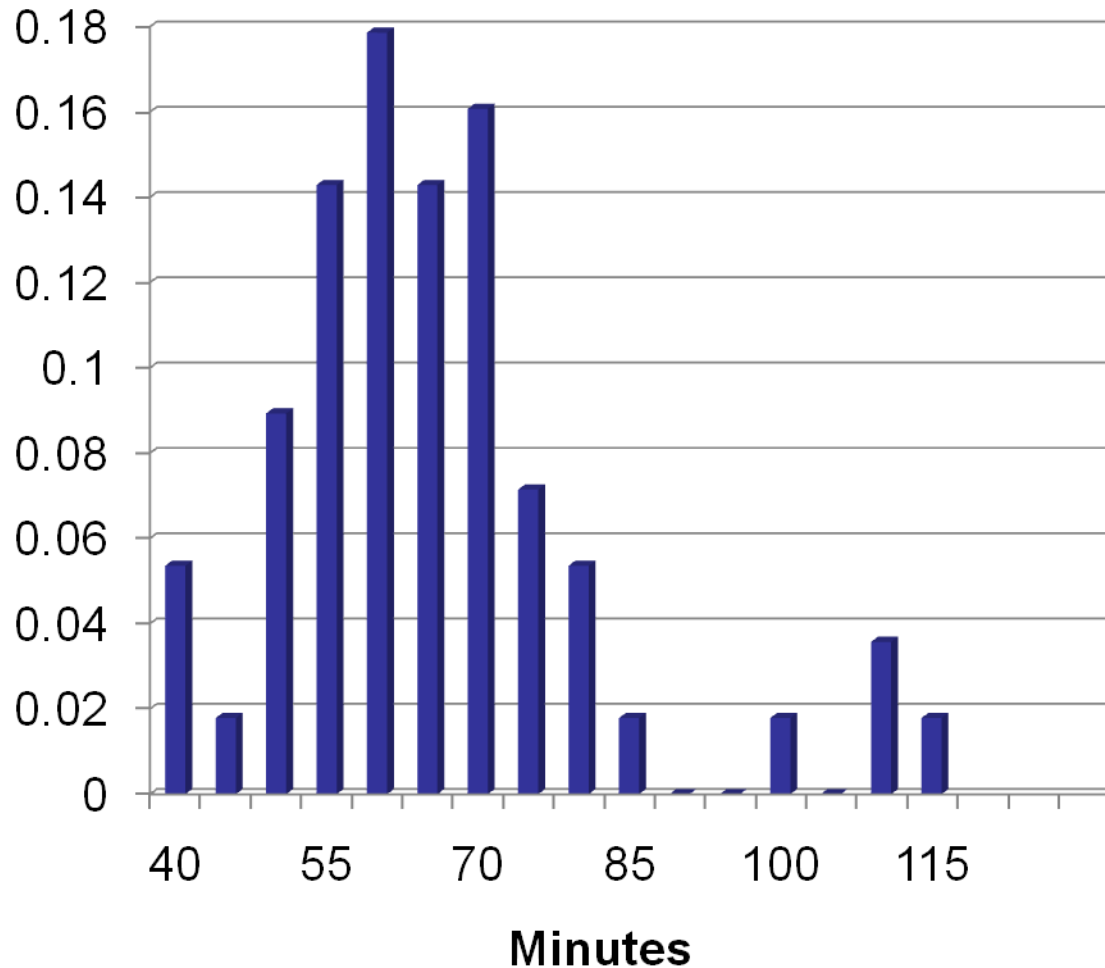
Examples of stochastic scheduling problems

- What is the optimal assignment of surgeries to operating rooms at a hospital?
- What is the optimal schedule of deliveries of raw material inventory to a manufacturer?
- What is the optimal arrival schedule of cargo ships to a port?

Complicating Factors

- High cost of customer waiting and server idling and a fixed time to complete activities
- Large number of activities to be coordinated in a constrained environment
- Uncertainty in the duration of customer service, and server availability
- Human behavior

A motivating example – surgery



Example 1

Single Server Scheduling

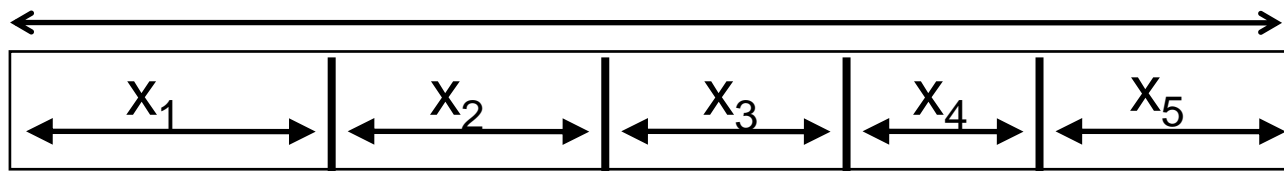
Single Server Scheduling Problem

For a single server, find the optimal time to allocate for each customer to minimize the cost of:

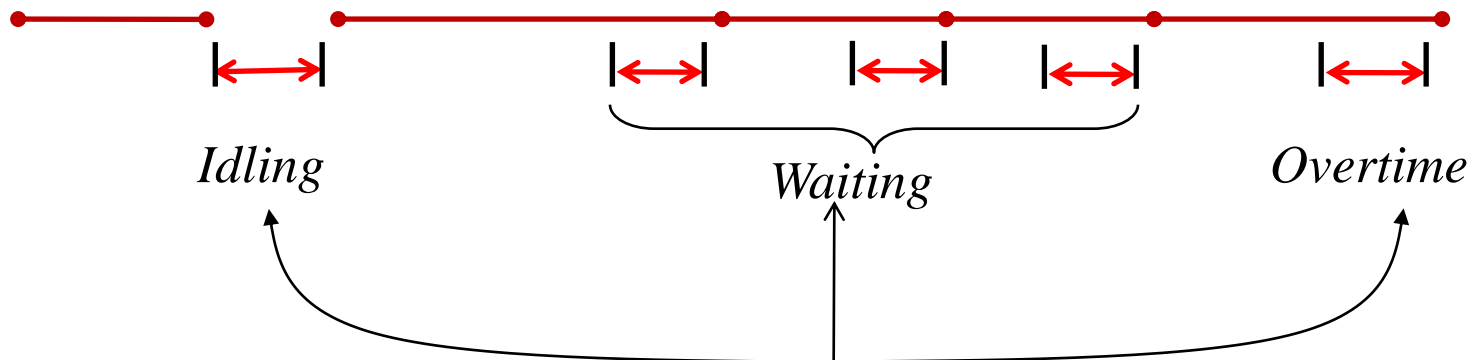
- Customer waiting
- Server idling
- Overtime

Single Server Scheduling

Planned Server Availability (e.g. 8 hours)



Example Scenario:



Goal: $\text{Min}\{\text{Idling} + \text{Waiting} + \text{Overtime}\}$

Stochastic Optimization Model

$$\min_x \left\{ \overbrace{\sum_{i=1}^n C_i^W E_Z[W_i]}^{\text{Cost of Waiting}} + \overbrace{\sum_{i=1}^n C_i^S E_Z[S_i]}^{\text{Cost of Idling}} + \overbrace{C^L E_Z[L]}^{\text{Cost of Overtime}} \right\}$$

Random
service time

Planned time
for service

$$W_i = \max(W_{i-1} + Z_{i-1} - x_{i-1}, 0)$$

$$S_i = \max(-W_{i-1} - Z_{i-1} + x_{i-1}, 0)$$

$$L = \max(W_n + Z_n + \sum x_i - T, 0)$$

Literature Review – Single Server

Queuing Analysis:

- Mercer (1960, 1973)
- Jansson (1966)
- Brahim and Worthington (1991)

Assumes steady state is reached,
i.i.d. service times

Heuristics:

- White and Pike (1964)
- Soriano (1966)
- Ho and Lau (1992)

No guarantee of optimal solution

Optimization:

- Weiss (1990) – 2 customer news vendor model
- Wang (1993) – Multiple customers with phase-type distribution property
- Denton and Gupta (2003) – General stochastic programming formulation

Reformulation as a Stochastic Program

$$\min_{\mathbf{x}} \{ E_Z [\sum_{i=2}^n c_i^W w_i + \sum_{i=2}^n c^S s_i + c^L l] \}$$

$$s.t. \quad w_2 \quad - s_2 \quad = Z_1 - x_1$$

$$-w_2 + w_3 \quad - s_3 \quad = Z_2 - x_2$$

.....

$$-w_n \quad - s_n + l - g = Z_n - d + \sum_{j=1}^{n-1} x_j$$

$$x_i \geq 0, w_i \geq 0, s_i \geq 0, i = 1, \dots, n, \quad l, g \geq 0$$

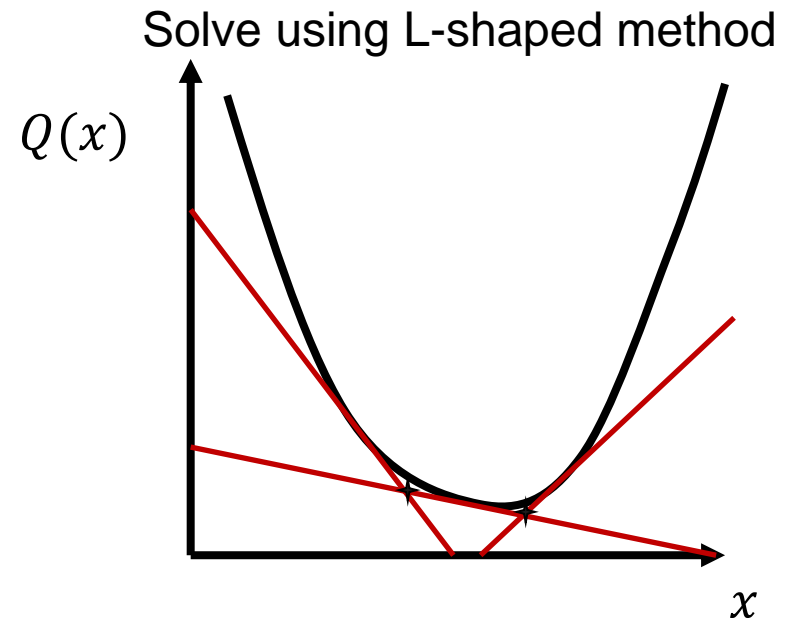
Two Stage Recourse Problem

Initial Decision (x) \rightarrow Uncertainty Resolved \rightarrow Recourse (y)

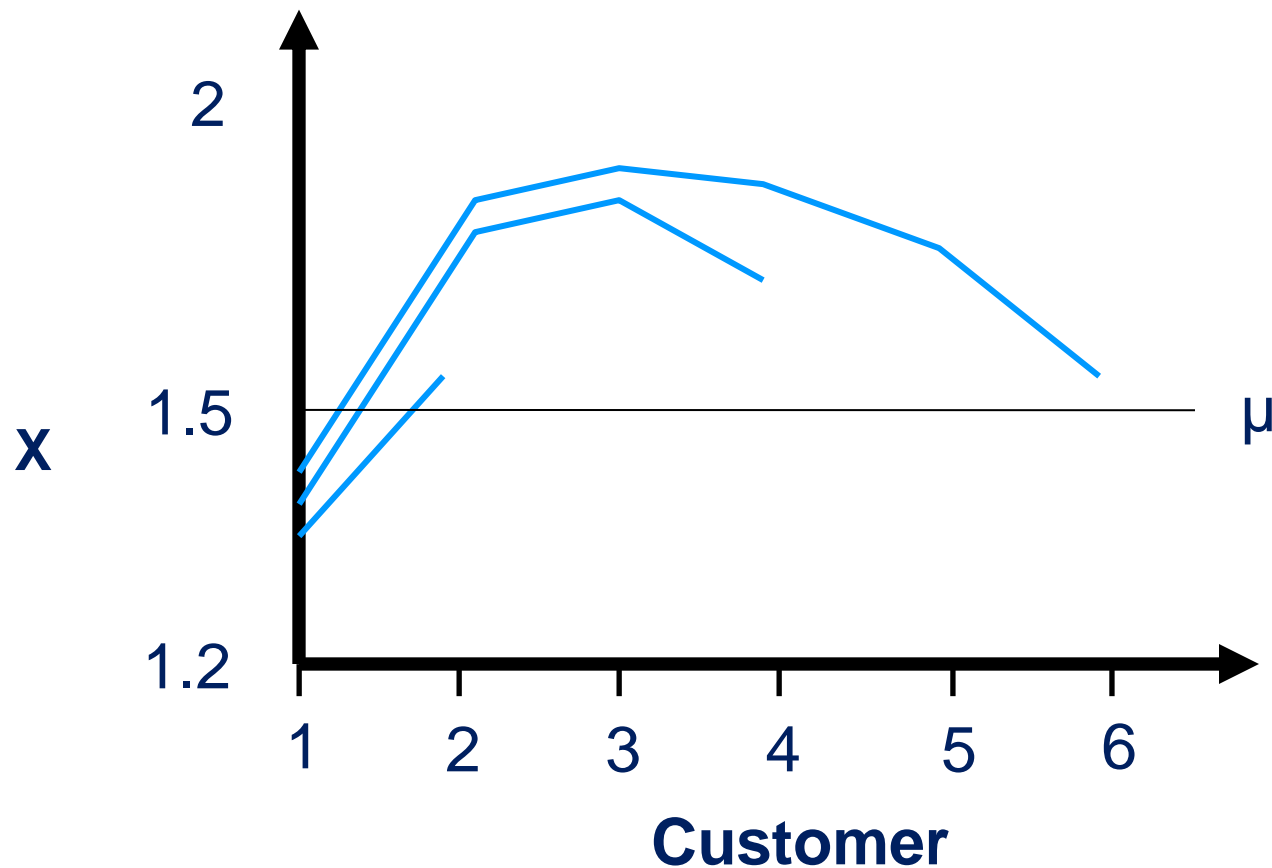
$$\min\{ Q(\mathbf{x}) = \sum_k^K p_k Q(\mathbf{x}, \mathbf{z}^k) \}$$

$$Q(\mathbf{x}, \mathbf{z}^k) = \min\{ \mathbf{c} \cdot \mathbf{y}^k \mid T \mathbf{x} + W \mathbf{y}^k = \mathbf{h}^k, \mathbf{y}^k \geq 0 \}$$

$$\left(\begin{array}{c} T \\ T \\ T \\ \vdots \\ T \end{array} \quad \begin{array}{c} W^1 \\ W^2 \\ W^3 \\ \dots \\ W^K \end{array} \right)$$



Example: $n = 3, 5, 7$ customers with i.i.d. service times $\sim U(1,2)$, $c^W = c^S$



Insights

- Simple heuristics often perform poorly
- The value of the stochastic solution (VSS) can be high
- Large instances of this problem can be solved very easily

1) Denton, B.T., Gupta, D., 2003, A Sequential Bounding Approach for Optimal Appointment Scheduling, *IIE Transactions*, 35, 1003-1016

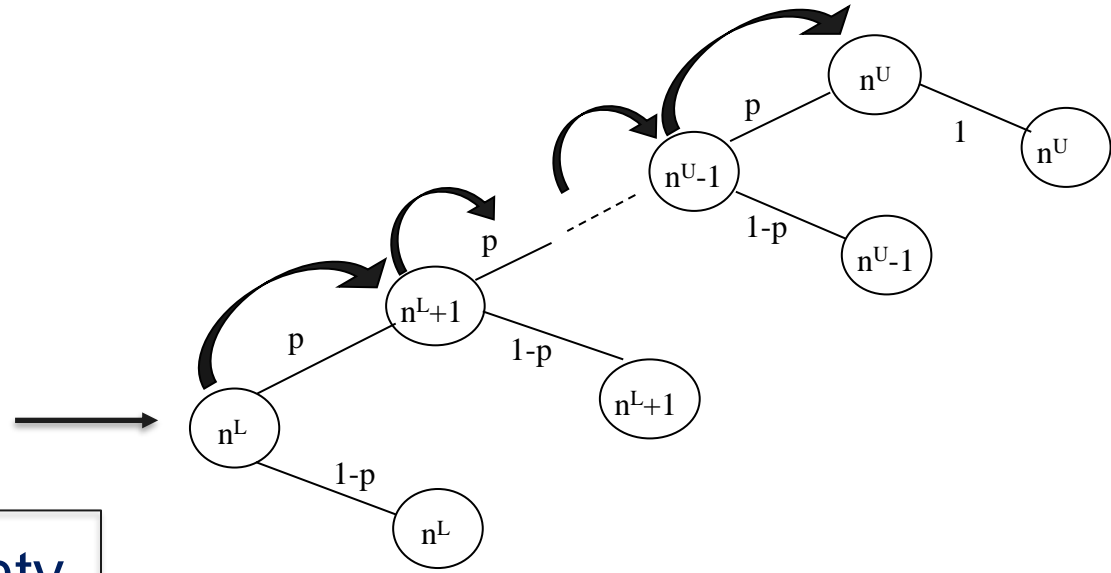
2) Denton, B.T., Viapiano, J, Vogl, A., 2007, Optimization of Surgery Sequencing and Scheduling Decisions Under Uncertainty, *Health Care Management Science*, 10(1), 13-24

There are many variations on this problem

- Customer No-shows
- Late arrivals

■ Dynamic scheduling

■ Endogenous uncertainty



Erdogan, S.A., Denton, B.T., “Dynamic Appointment Scheduling with Uncertain Demand,” *INFORMS Journal on Computing* 25(1), 116-132, 2013.

Erdogan, A, Denton, B.T., Gose, “On-line Appointment Sequencing and Scheduling,” *IIE Transactions*, 47, 1267-1286, 2015.

Zheng, Z., Denton, B.T., Xie, X., “Appointment Scheduling and the Effects of Customer Congestion on Service,” *IIE Transactions*; 51(10), 1075-1090, 2019

Example 2

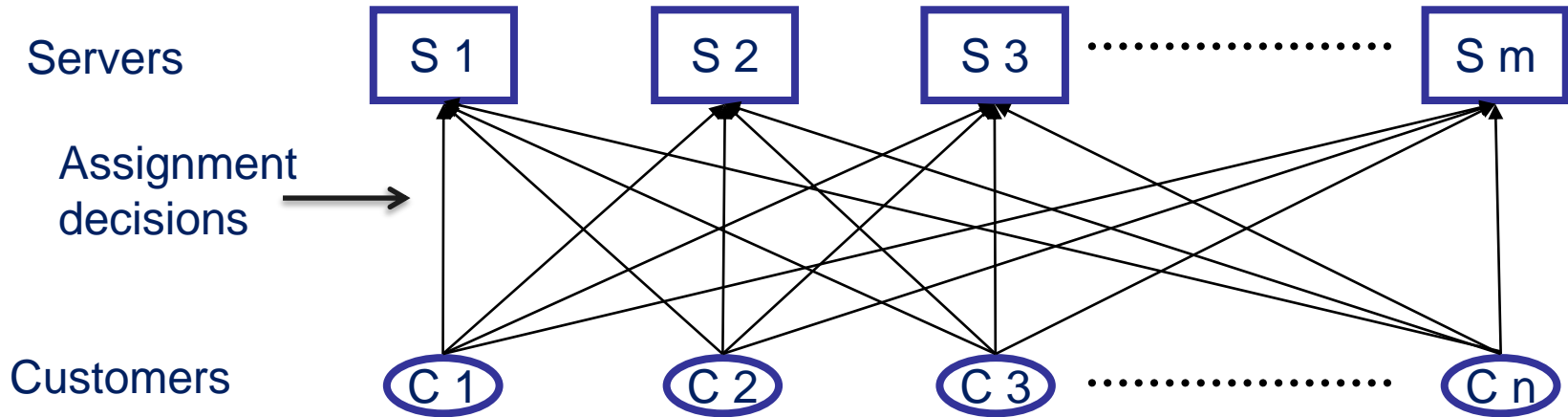
Multiple Server Job Allocation

Multi-Server Scheduling Problem

Given a set of customers (jobs) with uncertain duration to be scheduled on a certain day decide the following:

- How many servers to make available to complete all customer service
- Which server to assign to each customer

Multi-Server Scheduling Problem



Decisions:

- How many servers to have active each day?
- Which server to assign each job?

Extensible Bin-Packing Problem

$$x_i = \begin{cases} 1 & \text{if server } i \text{ active} \\ 0 & \text{otherwise} \end{cases}$$

$$y_{ij} = \begin{cases} 1 & \text{if customer } j \text{ assigned to server } i \\ 0 & \text{otherwise} \end{cases}$$

$$Z = \min \left\{ \sum_{i=1}^m c^f x_i + c^v o_i \right\}$$



Cost of servers + overtime

$$s. t. \quad y_{ij} \leq x_i \quad i = 1, \dots, m, j = 1, \dots, n$$



Customers only scheduled to active servers

$$\sum_{i=1}^m y_{ij} = 1 \quad j = 1, \dots, n$$



Every customer assigned to one server

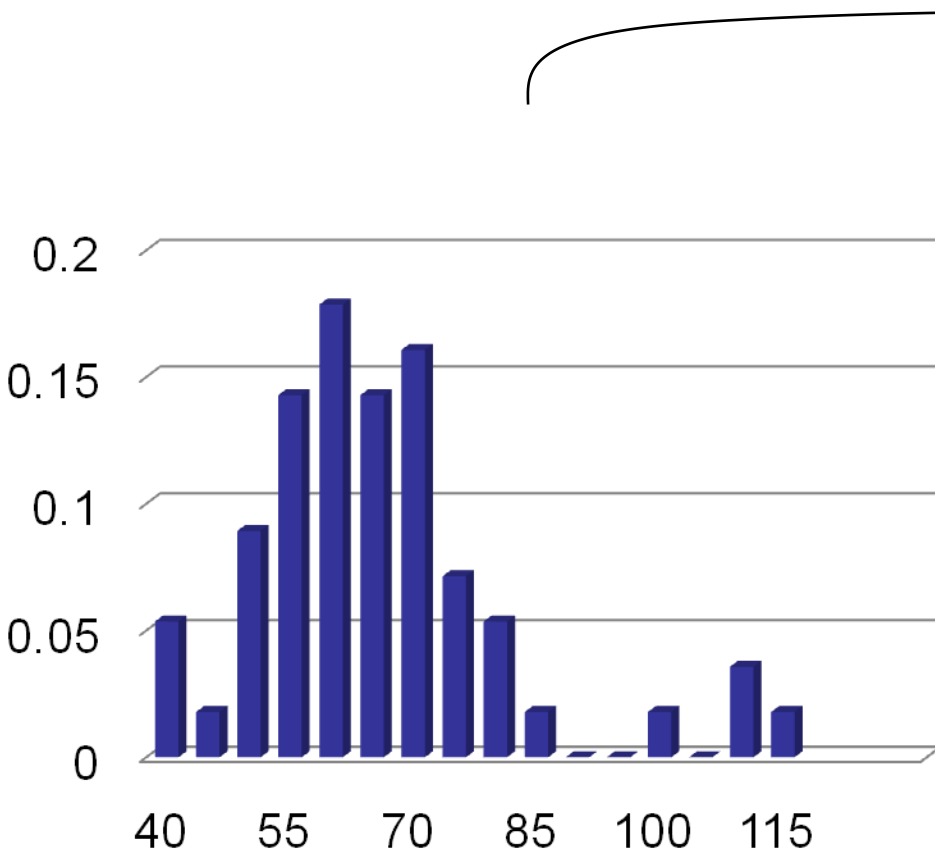
$$\sum_{j=1}^n d_j y_{ij} - o_i \leq T x_i \quad i = 1, \dots, m$$



Overtime if server goes past end of day of length T

$$y_{ij}, x_i \text{ binary}, \quad o_i \geq 0$$

Two-stage stochastic mixed integer program



$$Q(\mathbf{x}) = \min \left\{ \sum_{j=1}^m c^f x_j + c^v E_{\omega} [o_j(\omega)] \right\}$$

$$s.t. \quad y_{ij} \leq x_j \quad \forall (i, j)$$

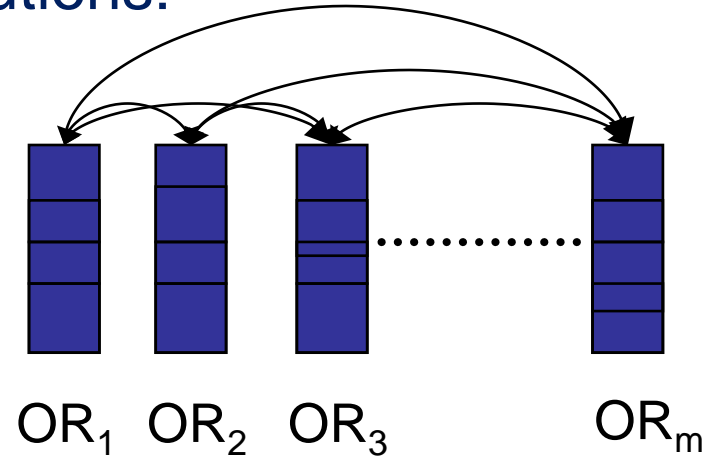
$$\sum_{j=1}^m y_{ij} = 1 \quad \forall (i)$$

$$\sum_{i=1}^n d_i(\omega) y_{ij} - o_j(\omega) \leq T x_j \quad \forall (i, j, \omega)$$

$$y_{ij}, x_j \in \{0, 1\}, \quad o_j(\omega) \geq 0, \forall \omega$$

Symmetry is a problem

There are $m!$ optimal solutions:



Adding the following anti-symmetry constraints reduces computation time:

$$\begin{aligned}x_1 &\geq x_2 \\x_2 &\geq x_3 \\&\vdots \\x_m &\geq x_{m-1}\end{aligned}$$

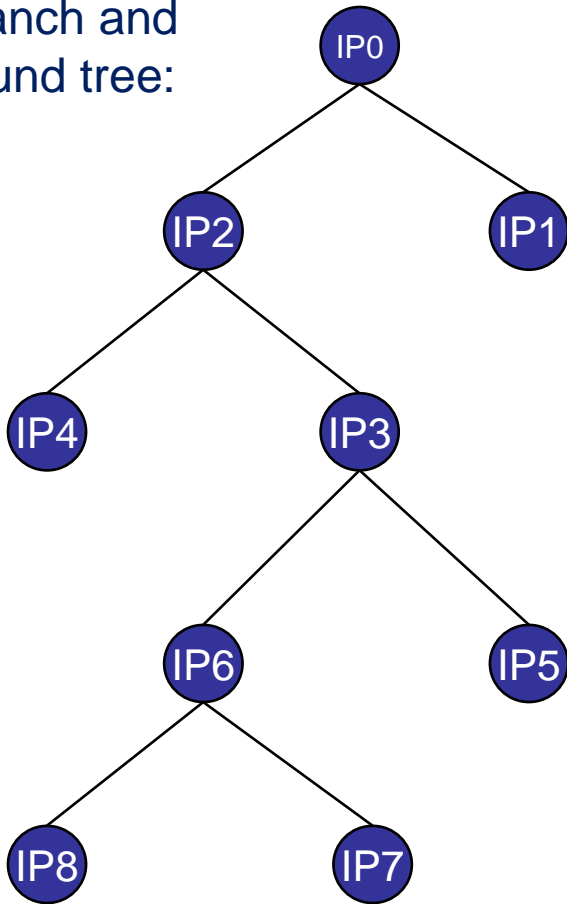
**Server
Ordering**

$$\begin{aligned}y_{11} &= 1 \\y_{21} + y_{22} &= 1 \\&\vdots \\ \sum_{j=1}^m y_{mj} &= 1\end{aligned}$$

**Customer
Assignment**

Integer L-Shaped Method

Branch and bound tree:



Master Problem:

$$Z = \min \left\{ \sum_{j=1}^m c^f x_j + \Theta \right\}$$

$$s. t. \quad y_{ij} \leq x_j \quad \forall (i, j)$$

$$\sum_{j=1}^m y_{ij} = 1 \quad \forall (i)$$

$$y_{ij}, x_j \in \{0, 1\}, \Theta \geq 0$$

(optimality cuts)

$$\Theta \geq E_{\omega}[\pi(h - Tx)]$$

Longest Processing Time First Heuristic

```
Sort customers in LPT order;  
 $m \leftarrow LB$  on number of servers;  
while( $o_j = 0, \forall j$ )  
    LPT( $m$ );  
     $m \leftarrow m + 1$ ;  
end  
Compute  $m^*$  with lowest total cost
```

Dell'Ollmo, Kellerer, Speranza, Tuza, *Information Processing Letters* (1998) – provides a 13/12 approximation algorithm for bin packing with a fixed number of extensible bins

Robust Formulation

Robust formulation seeks to minimize the worst case cost.

$$\begin{aligned}
 Z &= \min \left\{ \sum_{j=1}^m c^f x_j + F(x, y) \right\} \\
 \text{s. t. } & y_{ij} \leq x_j \quad \forall (i, j) \\
 & \sum_{j=1}^m y_{ij} = 1 \quad \forall (i) \\
 & y_{ij}, x_j \in \{0, 1\} \geq 0
 \end{aligned}$$

Worst case (adversary) problem

$$F(x, y) = \left\{ \begin{aligned}
 & \max_{\delta} \left\{ \sum_{j=1}^m \eta_j \right\} \\
 \text{s. t. } & \eta_j = c_j^v \max \left\{ 0, \sum_{i: y_{ij}=1} \delta_{ij} y_{ij} - dx_j \right\}, \quad \forall j \\
 & \sum_{(i,j): y_{ij}=1} \frac{\delta_{ij} - \underline{z}_i}{\bar{z}_i - \underline{z}_i} y_{ij} \leq \tau \quad \leftarrow \text{Uncertainty budget} \\
 & \underline{z}_i \leq \delta_{ij} \leq \bar{z}_i, \quad \forall (i, j): y_{ij} = 1
 \end{aligned} \right.$$

Results from sample test problems

Instance	1	2	3	4	5	6	7	8	9	10	Avg
LPT	.82	.97	.85	.93	.95	.85	.94	.97	.97	.92	.92
MV	.81	.95	.85	.92	.90	.86	.93	.89	.96	.86	.90
Robust	.93	.97	.97	.92	.89	.94	.92	.90	.97	.92	.92

Table 1: Cost of 0.5 hours overtime equals cost, c^f , of adding a server

Instance	1	2	3	4	5	6	7	8	9	10	Avg
LPT	1.0	1.0	1.0	1.0	1.0	.99	.99	.97	.99	1.0	.99
MV	1.0	1.0	1.0	1.0	.99	.99	.97	.97	.98	1.0	.99
Robust	.95	1.0	.95	.93	.94	.88	.97	.99	.96	.90	.95

Table 2: Cost of 2 hours overtime equals cost, c^f , of adding a server

LPT = longest processing time first heuristic, MV = mean value problem, Robust = solution to robust integer program. Results expressed as the ratio of optimal solution to solution generated by MV, LPT, Robust

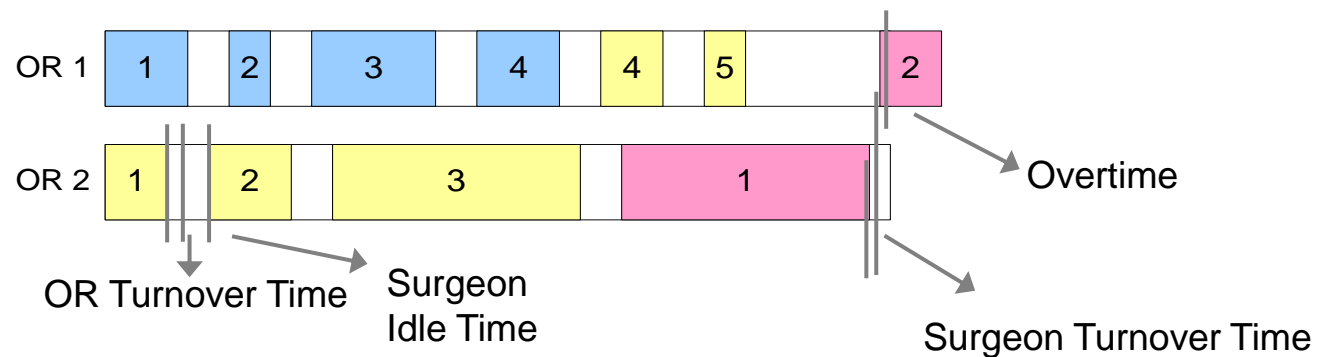
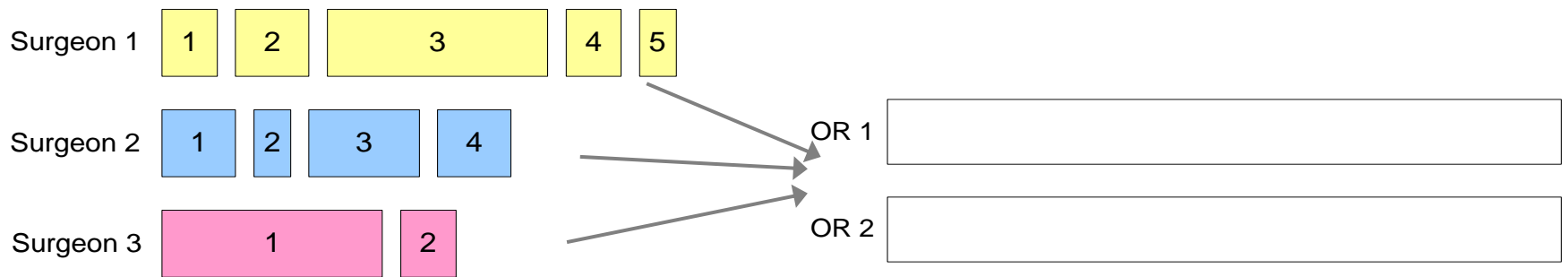
Insights

- LPT works well when overtime costs are low and it has a favorable *performance ratio*
- LPT is better (and much easier) than solving MV problem in most cases
- Robust IP is better than LPT when overtime costs are high

Denton, B.T., Miller, A., Balasubramanian, H., Huschka, T., 2010, Optimal Surgery Block Allocation Under Uncertainty, *Operations Research* 58(4), 802-816, 2010

Zheng, Z., Denton, B.T., Xie, X., “Branch-and-Price for Chance-Constrained Bin Packing,” *INFORMS Journal on Computing*; 32(3):547-564, 2020

Relaxing assumptions about assignment decisions leads to challenging problems



Batun, S., Denton, B.T., Huschka, T.R., Schaefer, A.J., The Benefit of Pooling Operating Rooms Under Uncertainty, *INFORMS Journal on Computing*, 23(2), 220-237, 2012.

Example 3

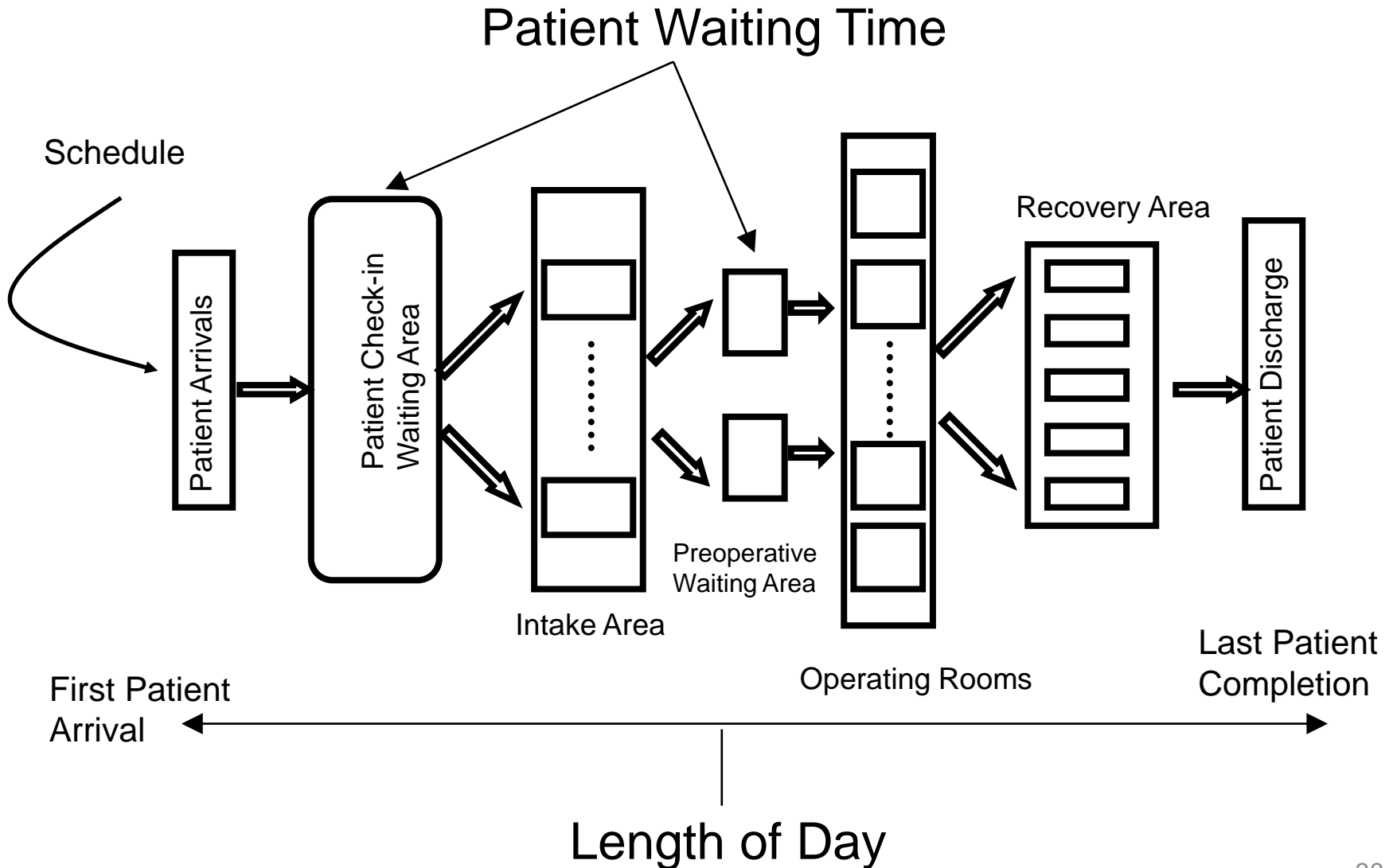
Multi-Stage Service System

Multi-Stage Server Scheduling Problem

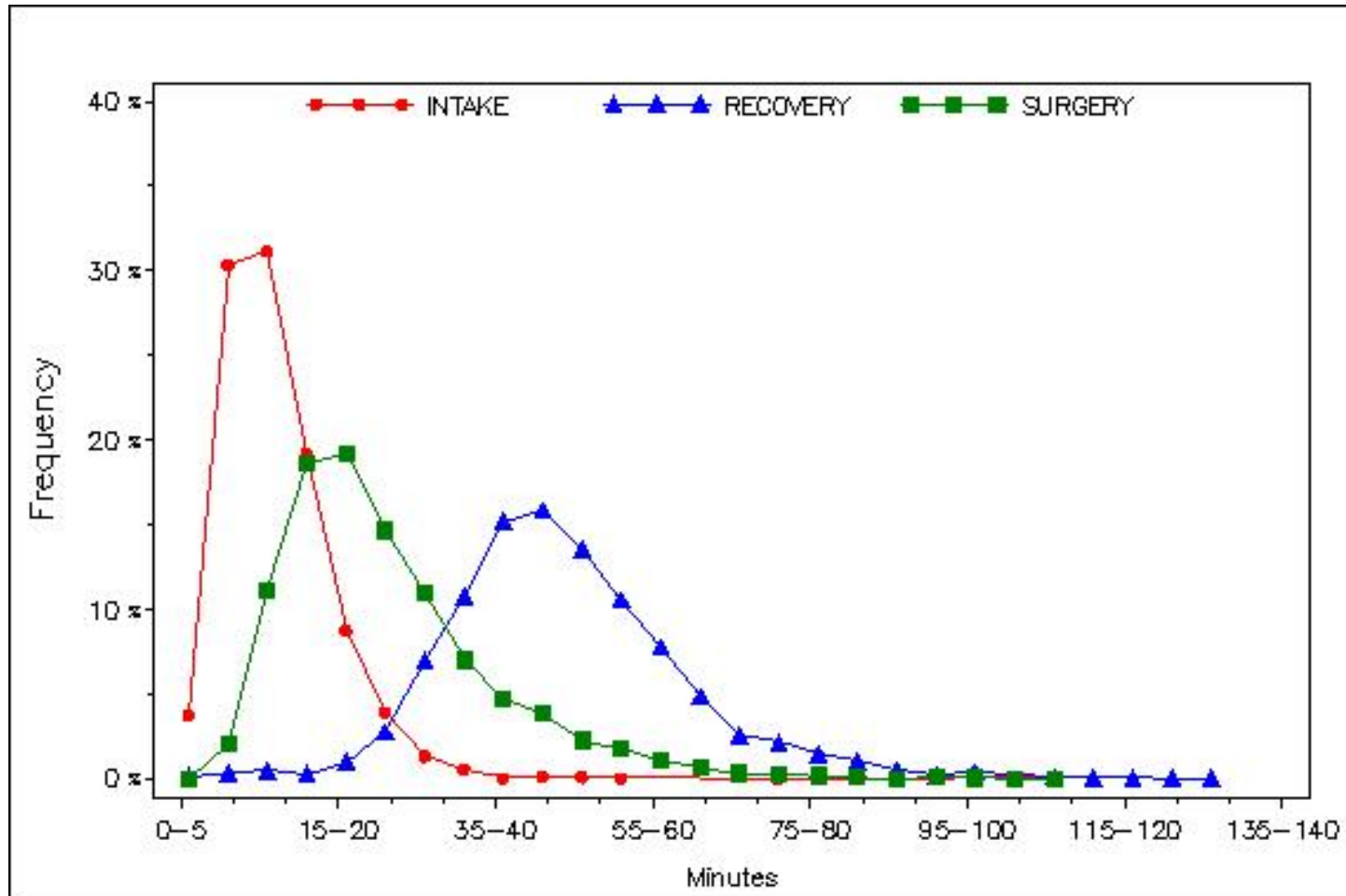
Find the Pareto optimal appointment times for patients having a procedure in an ambulatory surgery center to trade-off:

- Expected patient waiting
- Expected length of day

Context: Outpatient Procedure Centers



Intake, Procedure and Recovery Distributions



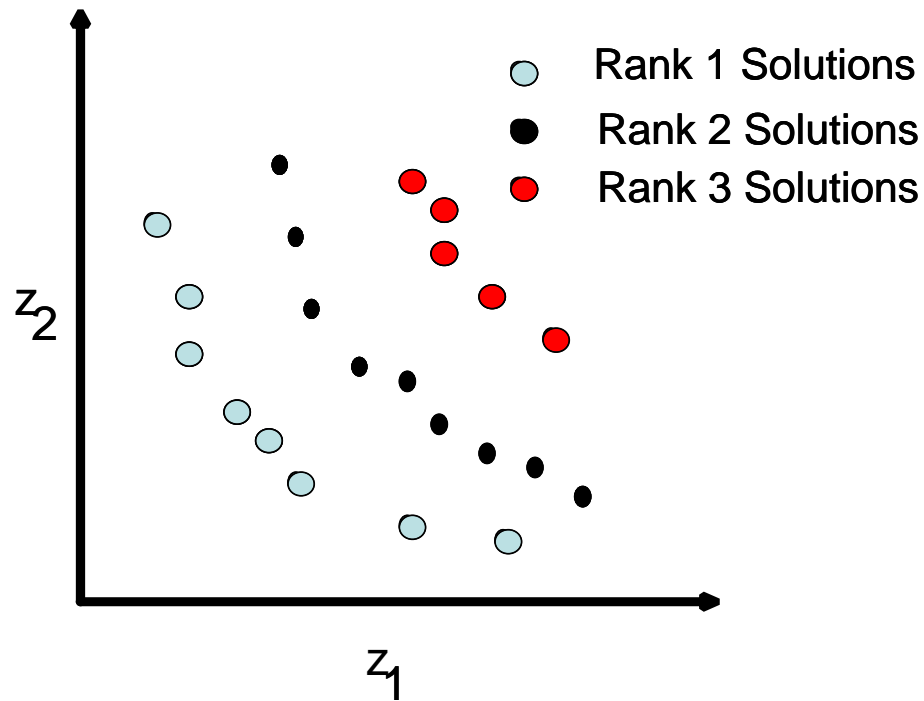
Simulation-optimization

Decision variables: scheduled arrival times to be assigned to n patients each day

Goal: Generate Pareto optimal schedules to understand tradeoffs between patient waiting and length of day

- Schedules generated using a genetic algorithm (GA)
- Non-dominated sorting used to identify the Pareto set and feedback into GA

The non-dominated sorting genetic algorithm (NSGA-II) of Deb *et al.*(2000):



Selection Procedure

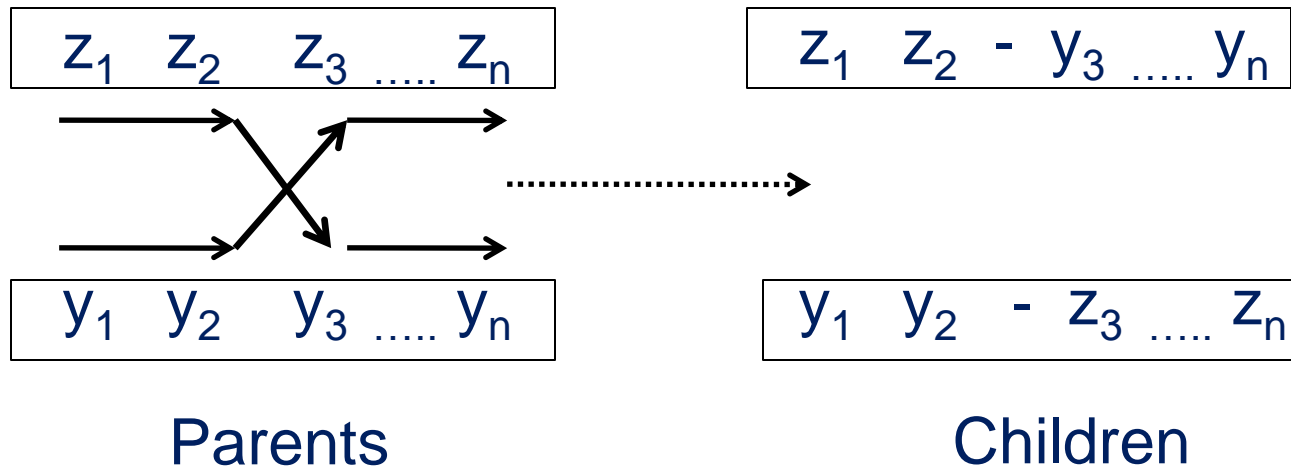
Sequential two stage indifference zone ranking and selection procedure of Rinott (1978) to compute the number of samples necessary to determine whether a solution i “dominates” j

Solution i “dominates” j if:

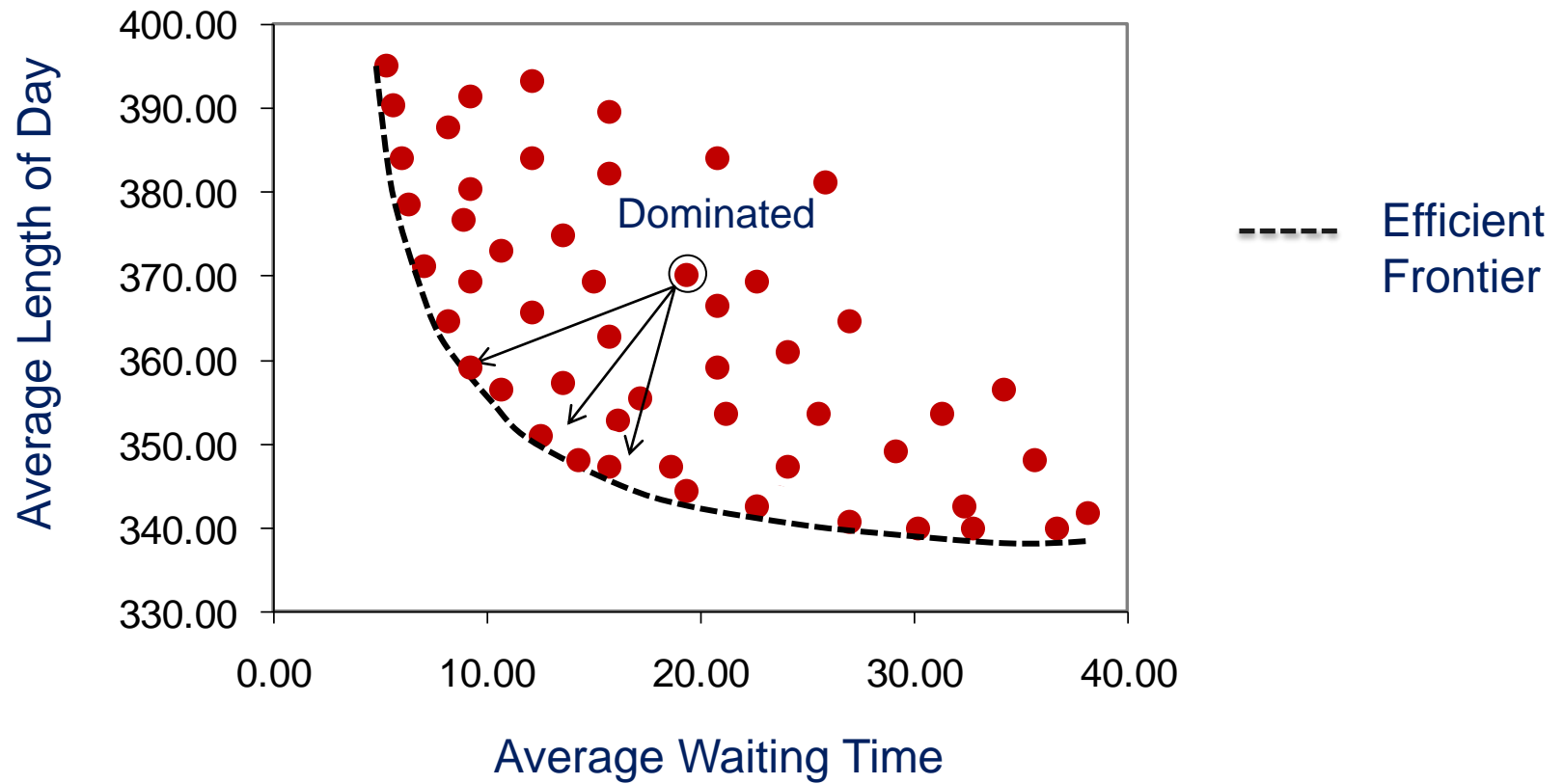
$$E[W_i] < E[W_j] \quad \text{and} \quad E[L_i] < E[L_j]$$

Genetic Algorithm

- Randomly generated initial population of schedules
- Selection based on 1) ranks and 2) crowding distance
- Mutation
- Single point crossover:



Schedule Optimization



Insights

- A simple simulation-optimization approach provides significant improvement to schedules used in practice
- Substantial reduction in average waiting time is possible with a very limited increase in average length of day

Gul, S., Denton, B.T., Fowler, J., 2011 Bi-Criteria Scheduling of Surgical Services for an Outpatient Procedure Center, *Production and Operations Management*, 20(3), 406-417

Key Takeaways

- Modeling uncertainty often matters!
- Stochastic scheduling problems can be hard, but a special structure often exists to exploit.
- Stochastic optimization is a powerful tool for scheduling in many contexts



Acknowledgements

Hari Balasubramanian (University of Massachusetts)

Sakine Batun (METU, Turkey)

Maya Bam (General Motors)

Bjorn Berg (University of Minnesota)

Ayca Erdogan (San Jose State University)

Serhat Gul (TED University, Turkey)

Todd Huschka (Mayo Clinic)

Andrew Miller (FedEx)

Heidi Nelson (Mayo Clinic)

Andrew Schaefer (Rice University)

Zheng Zhang (Zhejiang University)

Some of this work was funded in part by grants from the *Service Enterprise Systems* program at the *National Science Foundation*

Thank You!

Brian Denton
University of Michigan

btdenton@umich.edu



These slides and the papers cited can be found at my website:

<http://umich.edu/~btdenton>