

Tools for Staffing and Scheduling in Call Centers: Survey of work in our lab

Pierre L'Ecuyer

with Thanos Avramidis, Eric Buist, Tolga Cezik, Wyeon Chan, Nabil Channouf
Canada Research Chair in Stochastic Simulation and Optimization, U. Montréal

Sponsored by Bell Canada

Tools for Staffing and Scheduling in Call Centers: Survey of work in our lab

Pierre L'Ecuyer

with Thanos Avramidis, Eric Buist, Tolga Cezik, Wyeon Chan, Nabil Channouf
Canada Research Chair in Stochastic Simulation and Optimization, U. Montréal

Sponsored by Bell Canada

- General problem statement;

Tools for Staffing and Scheduling in Call Centers: Survey of work in our lab

Pierre L'Ecuyer

with Thanos Avramidis, Eric Buist, Tolga Cezik, Wyeon Chan, Nabil Channouf
Canada Research Chair in Stochastic Simulation and Optimization, U. Montréal

Sponsored by Bell Canada

- General problem statement;
- Simulation tools and improving simulation efficiency;
Queueing approximations vs simulation;

Tools for Staffing and Scheduling in Call Centers:

Survey of work in our lab

Pierre L'Ecuyer

with Thanos Avramidis, Eric Buist, Tolga Cezik, Wyeon Chan, Nabil Channouf

Canada Research Chair in Stochastic Simulation and Optimization, U. Montréal

Sponsored by Bell Canada

- General problem statement;
- Simulation tools and improving simulation efficiency;
Queueing approximations vs simulation;
- Building realistic models; Computer games vs tools for decision making;

Tools for Staffing and Scheduling in Call Centers:

Survey of work in our lab

Pierre L'Ecuyer

with Thanos Avramidis, Eric Buist, Tolga Cezik, Wyeon Chan, Nabil Channouf
Canada Research Chair in Stochastic Simulation and Optimization, U. Montréal

Sponsored by Bell Canada

- General problem statement;
- Simulation tools and improving simulation efficiency;
Queueing approximations vs simulation;
- Building realistic models; Computer games vs tools for decision making;
- Optimization of staffing, scheduling, call routing, priorities, outbound call, etc.

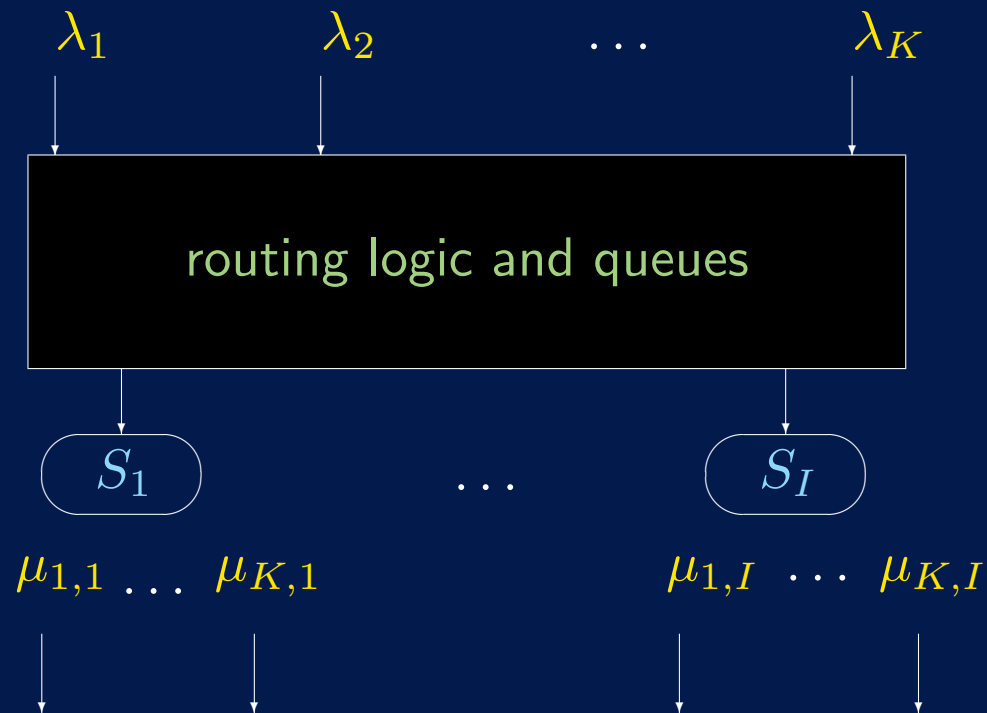
Article are on my web page: type Pierre L'Ecuyer in Google.

Example: A Call Center with Multiple Call Types

K call types. Depends on required technical skill, language, importance, etc.
 I agent types (or skill groups). Each has skills to handle certain call types.
Some agents may be better than others for a given call type.

Example: A Call Center with Multiple Call Types

K call types. Depends on required technical skill, language, importance, etc.
 I agent types (or skill groups). Each has skills to handle certain call types.
Some agents may be better than others for a given call type.



Call arrivals: follow some stochastic process.

Call arrivals: follow some stochastic process.

Abandonments: callers may abandon after a random patience time.

Call arrivals: follow some stochastic process.

Abandonments: callers may abandon after a random patience time.

If each agent has a single skill: K single queues in parallel.

Call arrivals: follow some stochastic process.

Abandonments: callers may abandon after a random patience time.

If each agent has a single skill: K single queues in parallel.

If each agent has all skills: one single queue.

More efficient (less wait, fewer abandonments), but more costly.

Call arrivals: follow some stochastic process.

Abandonments: callers may abandon after a random patience time.

If each agent has a *single skill*: K single queues in parallel.

If each agent has *all skills*: one single queue.

More efficient (less wait, fewer abandonments), but more costly.

Each additional skill increases the cost of an agent.

Can also decrease its speed!

Call arrivals: follow some stochastic process.

Abandonments: callers may abandon after a random patience time.

If each agent has a *single skill*: K single queues in parallel.

If each agent has *all skills*: one single queue.

More efficient (less wait, fewer abandonments), but more costly.

Each additional skill increases the cost of an agent.

Can also decrease its speed!

For well-balanced systems, one or two skills per agent often gives a performance almost as good as all skills for all agents (e.g., Wallace and Whitt 2004).

Performance measures

Example: Total cost of agents.

Performance measures

Example: Total cost of agents.

Constraints

Service level (**SL**): fraction of calls that wait less than an acceptable wait time ℓ (typically 20 to 30 seconds). Often measured and controlled separately by time period (hour, day, ...), by call type, + aggregated.

Performance measures

Example: Total cost of agents.

Constraints

Service level (**SL**): fraction of calls that wait less than an acceptable wait time ℓ (typically 20 to 30 seconds). Often measured and controlled separately by time period (hour, day, ...), by call type, + aggregated.

Caveat: could cheat and never serve calls that already waited too much.

Performance measures

Example: Total cost of agents.

Constraints

Service level (SL): fraction of calls that wait less than an acceptable wait time ℓ (typically 20 to 30 seconds). Often measured and controlled separately by time period (hour, day, ...), by call type, + aggregated.

Caveat: could cheat and never serve calls that already waited too much.

Conditional SL: $\mathbb{E}[\text{Wait} \mid \text{Wait} > \ell]$.

Performance measures

Example: Total cost of agents.

Constraints

Service level (SL): fraction of calls that wait less than an acceptable wait time ℓ (typically 20 to 30 seconds). Often measured and controlled separately by time period (hour, day, ...), by call type, + aggregated.

Caveat: could cheat and never serve calls that already waited too much.

Conditional SL: $\mathbb{E}[\text{Wait} \mid \text{Wait} > \ell]$.

Abandonment ratio: fraction of calls that abandon.

Performance measures

Example: Total cost of agents.

Constraints

Service level (SL): fraction of calls that wait less than an acceptable wait time ℓ (typically 20 to 30 seconds). Often measured and controlled separately by time period (hour, day, ...), by call type, + aggregated.

Caveat: could cheat and never serve calls that already waited too much.

Conditional SL: $\mathbb{E}[\text{Wait} \mid \text{Wait} > \ell]$.

Abandonment ratio: fraction of calls that abandon.

Average waiting time?

Skill-based routing (SBR) strategies:

Rules that control in real time the agent-to-call and call-to-agent assignments.
Act as controlling device for service levels, across call types and globally.

Skill-based routing (SBR) strategies:

Rules that control in real time the agent-to-call and call-to-agent assignments.
Act as controlling device for service levels, across call types and globally.

Dynamic routing: Decision may depend on the entire state of the system.
An optimal policy is generally too complicated and hard to implement.

Skill-based routing (SBR) strategies:

Rules that control in real time the agent-to-call and call-to-agent assignments. Act as controlling device for service levels, across call types and globally.

Dynamic routing: Decision may depend on the entire state of the system. An optimal policy is generally too complicated and hard to implement.

Static routing: Each call type has an ordered list of agent types. If all are busy, the call joins a queue (usually one queue per call).

Skill-based routing (SBR) strategies:

Rules that control in real time the agent-to-call and call-to-agent assignments. Act as controlling device for service levels, across call types and globally.

Dynamic routing: Decision may depend on the entire state of the system. An optimal policy is generally too complicated and hard to implement.

Static routing: Each call type has an ordered list of agent types. If all are busy, the call joins a queue (usually one queue per call).

Each agent type may also have an ordered list of queues (call types) for when it becomes available (priorities).

Skill-based routing (SBR) strategies:

Rules that control in real time the agent-to-call and call-to-agent assignments. Act as controlling device for service levels, across call types and globally.

Dynamic routing: Decision may depend on the entire state of the system. An optimal policy is generally too complicated and hard to implement.

Static routing: Each call type has an ordered list of agent types. If all are busy, the call joins a queue (usually one queue per call).

Each agent type may also have an ordered list of queues (call types) for when it becomes available (priorities).

Other: Could weight the calls, use state-dependent thresholds, etc.

Staffing problem: Divide the day into periods (e.g, half hours) and determine how many agents of each type to have in the center for each period in order to meet the performance constraints, at minimal cost.

Staffing problem: Divide the day into periods (e.g, half hours) and determine how many agents of each type to have in the center for each period in order to meet the performance constraints, at minimal cost.

The SL in one period can depend on the number of agents in other periods!

Staffing problem: Divide the day into periods (e.g, half hours) and determine how many agents of each type to have in the center for each period in order to meet the performance constraints, at minimal cost.

The SL in one period can depend on the number of agents in other periods!

The SL constraints can be per customer type, per period, or aggregated.

Staffing problem: Divide the day into periods (e.g, half hours) and determine how many agents of each type to have in the center for each period in order to meet the performance constraints, at minimal cost.

The SL in one period can depend on the number of agents in other periods!

The SL constraints can be per customer type, per period, or aggregated.

Caveat: not possible to match exactly the optimal staffing by scheduling a set of agents whose working shifts are admissible (i.e., satisfy the constraints determined by Union rules, etc.).

Staffing problem: Divide the day into periods (e.g, half hours) and determine how many agents of each type to have in the center for each period in order to meet the performance constraints, at minimal cost.

The SL in one period can depend on the number of agents in other periods!

The SL constraints can be per customer type, per period, or aggregated.

Caveat: not possible to match exactly the optimal staffing by scheduling a set of agents whose working shifts are admissible (i.e., satisfy the constraints determined by Union rules, etc.).

Scheduling problem: Determine a set of agents, each with its working shift for the day, so that the performance constraints are met, at minimal cost.

Staffing problem: Divide the day into periods (e.g, half hours) and determine how many agents of each type to have in the center for each period in order to meet the performance constraints, at minimal cost.

The SL in one period can depend on the number of agents in other periods!

The SL constraints can be per customer type, per period, or aggregated.

Caveat: not possible to match exactly the optimal staffing by scheduling a set of agents whose working shifts are admissible (i.e., satisfy the constraints determined by Union rules, etc.).

Scheduling problem: Determine a set of agents, each with its working shift for the day, so that the performance constraints are met, at minimal cost.

Scheduling and rostering problem:

In practice, we do not have an infinite supply of each agent type!

For a given set of agents and a given set of admissible shifts, assign a shift to each agent, to meet the performance constraints at minimal cost.

Blend systems: inbound and outbound calls

A predictive dialer makes outbound call, trying to reach customers when inbound traffic is low.

Blend systems: inbound and outbound calls

A predictive dialer makes outbound call, trying to reach customers when inbound traffic is low.

A right party connect: when the outbound contact is successful (good).

Blend systems: inbound and outbound calls

A predictive dialer makes outbound call, trying to reach customers when inbound traffic is low.

A right party connect: when the outbound contact is successful (good).

A mismatch: when the successful contact cannot be served immediately (bad).

Blend systems: inbound and outbound calls

A **predictive dialer** makes outbound call, trying to reach customers when inbound traffic is low.

A **right party connect**: when the outbound contact is successful (good).

A **mismatch**: when the successful contact cannot be served immediately (bad).

Possible **additional constraints**:

- Lower bound on the (expected) volume of outbound calls per day.
- Upper bound on the (expected) fraction of mismatches.

Blend systems: inbound and outbound calls

A predictive dialer makes outbound call, trying to reach customers when inbound traffic is low.

A right party connect: when the outbound contact is successful (good).

A mismatch: when the successful contact cannot be served immediately (bad).

Possible **additional constraints**:

- Lower bound on the (expected) volume of outbound calls per day.
- Upper bound on the (expected) fraction of mismatches.

Other types of recourses:

Ability to get new agents on short advice (perhaps on workstation at home)?

Blend systems: inbound and outbound calls

A predictive dialer makes outbound call, trying to reach customers when inbound traffic is low.

A right party connect: when the outbound contact is successful (good).

A mismatch: when the successful contact cannot be served immediately (bad).

Possible **additional constraints**:

- Lower bound on the (expected) volume of outbound calls per day.
- Upper bound on the (expected) fraction of mismatches.

Other types of recourses:

Ability to get new agents on short advice (perhaps on workstation at home)?

Controlling the arrival rate, e.g., by giving information on current waiting times?

Blend systems: inbound and outbound calls

A predictive dialer makes outbound call, trying to reach customers when inbound traffic is low.

A right party connect: when the outbound contact is successful (good).

A mismatch: when the successful contact cannot be served immediately (bad).

Possible **additional constraints**:

- Lower bound on the (expected) volume of outbound calls per day.
- Upper bound on the (expected) fraction of mismatches.

Other types of recourses:

Ability to get new agents on short advice (perhaps on workstation at home)?

Controlling the arrival rate, e.g., by giving information on current waiting times?

Other ways?

Simulation tools

Paradigm: a computer game that reproduces exactly the behavior of the call center relevant to the decisions we want to make.

Simulation tools

Paradigm: a computer game that reproduces exactly the behavior of the call center relevant to the decisions we want to make.

Then we can try everything we want without paying the price!

Simulation tools

Paradigm: a computer game that reproduces exactly the behavior of the call center relevant to the decisions we want to make.

Then we can try everything we want without paying the price!

Available software to do this?

Arena Contact Center Edition, ccProphet, ...

Simulation tools

Paradigm: a computer game that reproduces exactly the behavior of the call center relevant to the decisions we want to make.

Then we can try everything we want without paying the price!

Available software to do this?

Arena Contact Center Edition, ccProphet, ...

Plus: Graphical user interface and animation.

Simulation tools

Paradigm: a computer game that reproduces exactly the behavior of the call center relevant to the decisions we want to make.

Then we can try everything we want without paying the price!

Available software to do this?

Arena Contact Center Edition, ccProphet, ...

Plus: Graphical user interface and animation.

Minus: Expensive, slow, not so flexible.

Simulation tools

Paradigm: a computer game that reproduces exactly the behavior of the call center relevant to the decisions we want to make.

Then we can try everything we want without paying the price!

Available software to do this?

Arena Contact Center Edition, ccProphet, ...

Plus: Graphical user interface and animation.

Minus: Expensive, slow, not so flexible.

So we developed our own tool: **ContactCenters**, a Java library built over SSJ.

ContactCenters

Library based on well-supported modern programming language Java.

Strong expressive power. GUI-independent.

Interoperability with other software (statistics, optimization, databases, etc).

ContactCenters

Library based on well-supported modern programming language Java.

Strong expressive power. GUI-independent.

Interoperability with other software (statistics, optimization, databases, etc).

Fast: about 30 times faster than Arena in our experiments.

ContactCenters

Library based on well-supported modern programming language Java.

Strong expressive power. GUI-independent.

Interoperability with other software (statistics, optimization, databases, etc).

Fast: about 30 times faster than Arena in our experiments.

Provides building blocks to simulate all types of call centers.

Supports several types of contacts, multiskill, blend, arbitrary dialing and routing policies, various types of arrival processes, etc.

ContactCenters

Library based on well-supported modern programming language Java.

Strong expressive power. GUI-independent.

Interoperability with other software (statistics, optimization, databases, etc).

Fast: about 30 times faster than Arena in our experiments.

Provides building blocks to simulate all types of call centers.

Supports several types of contacts, multiskill, blend, arbitrary dialing and routing policies, various types of arrival processes, etc.

Support for variance reduction.

ContactCenters

Library based on well-supported modern programming language Java.

Strong expressive power. GUI-independent.

Interoperability with other software (statistics, optimization, databases, etc).

Fast: about 30 times faster than Arena in our experiments.

Provides building blocks to simulate all types of call centers.

Supports several types of contacts, multiskill, blend, arbitrary dialing and routing policies, various types of arrival processes, etc.

Support for variance reduction.

Pre-compiled generic models.

Building realistic models

We need lots of detailed data about the call center operations: arrivals, service times, breaks, hidden rules and practice, etc.

Building realistic models

We need lots of detailed data about the call center operations: arrivals, service times, breaks, hidden rules and practice, etc.

With partial and aggregate information, we can only build an “approximate” model, which may not be realistic.

Building realistic models

We need lots of detailed data about the call center operations: arrivals, service times, breaks, hidden rules and practice, etc.

With partial and aggregate information, we can only build an “approximate” model, which may not be realistic.

This is often a problem!

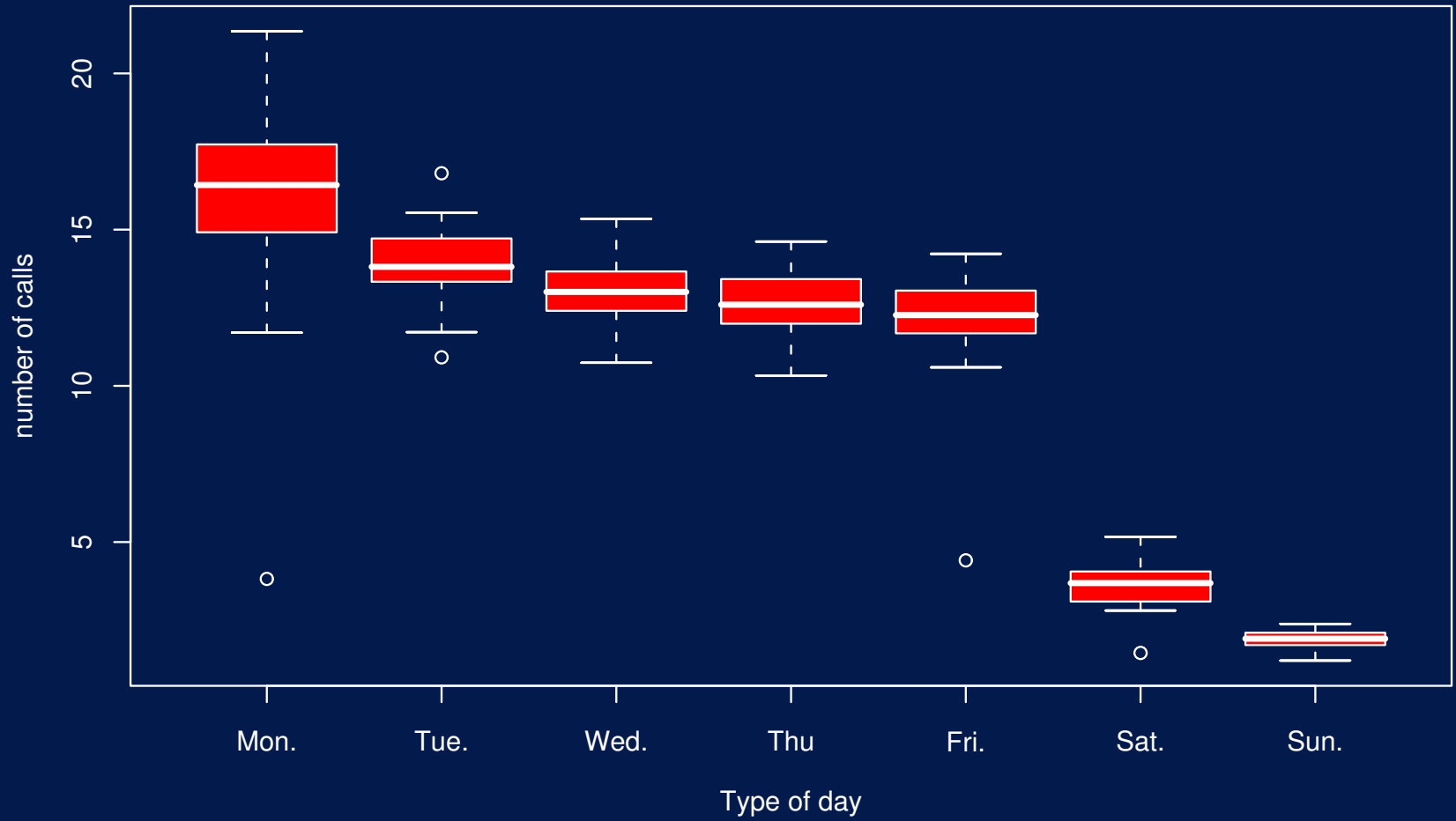
Building realistic models

We need lots of detailed data about the call center operations: arrivals, service times, breaks, hidden rules and practice, etc.

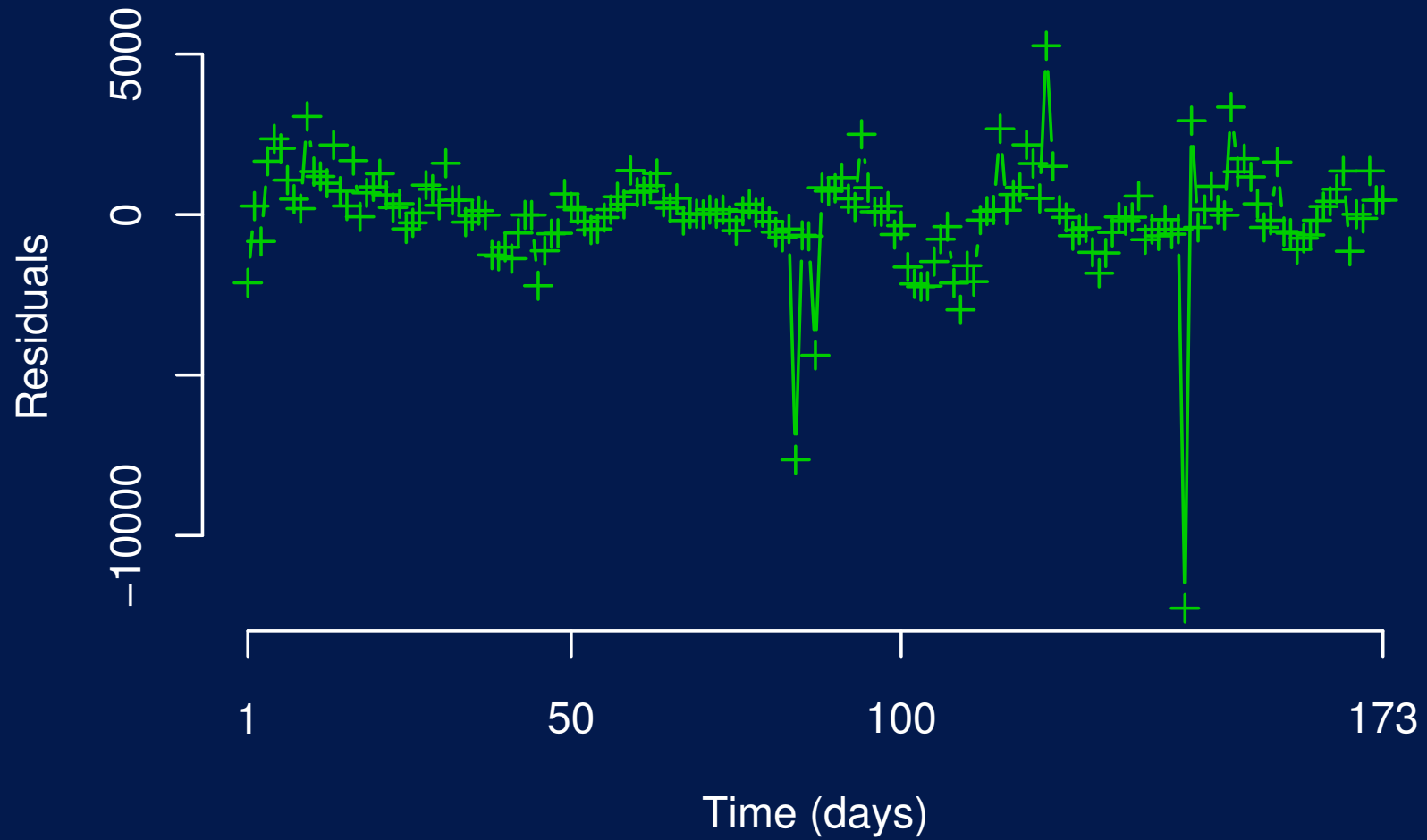
With partial and aggregate information, we can only build an “approximate” model, which may not be realistic.

This is often a problem!

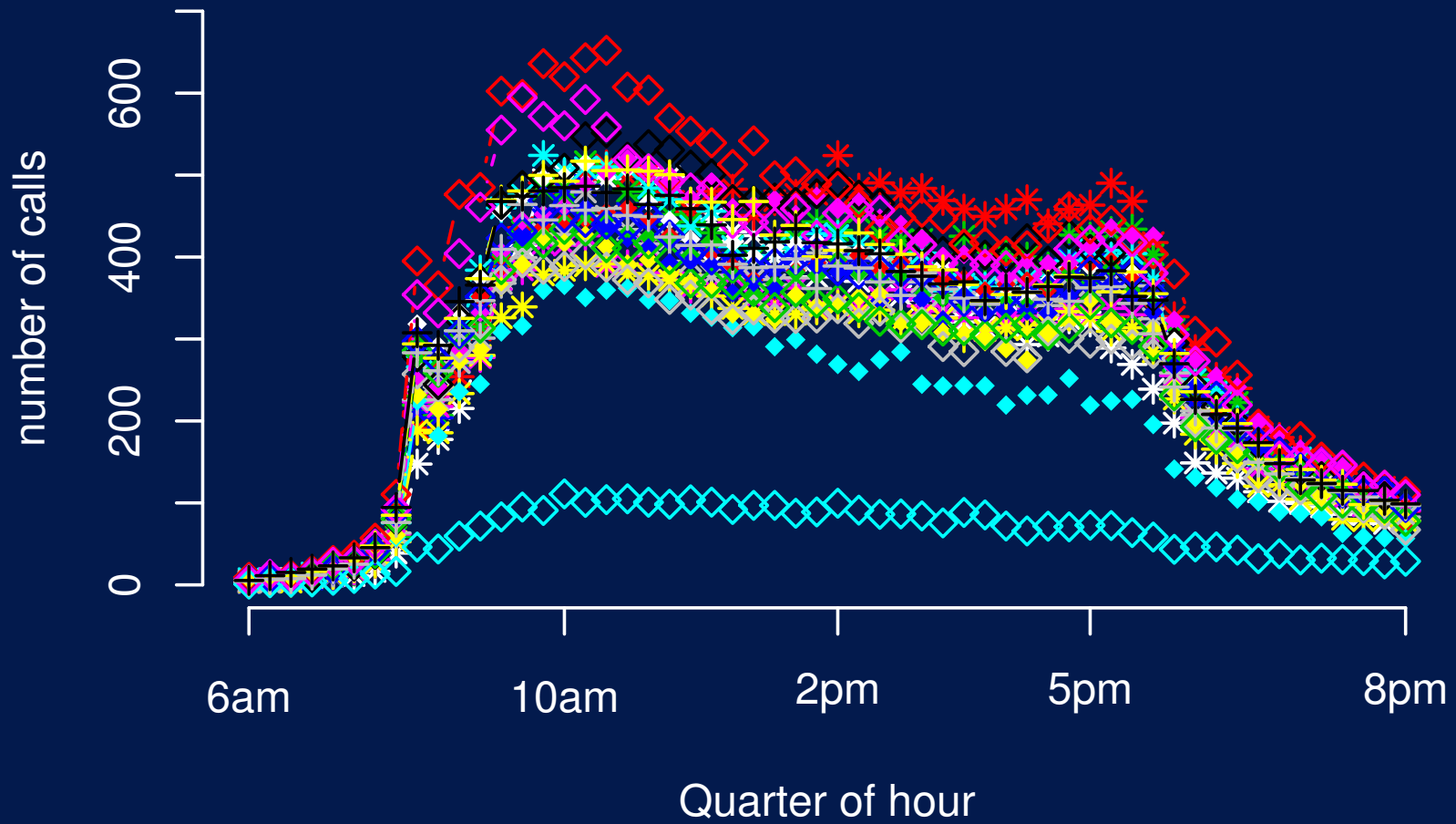
Examples: service times are not exponential; arrivals are not Poisson and stationary, perhaps breaks are not exactly as planned, perhaps agents are not available exactly as planned, etc.



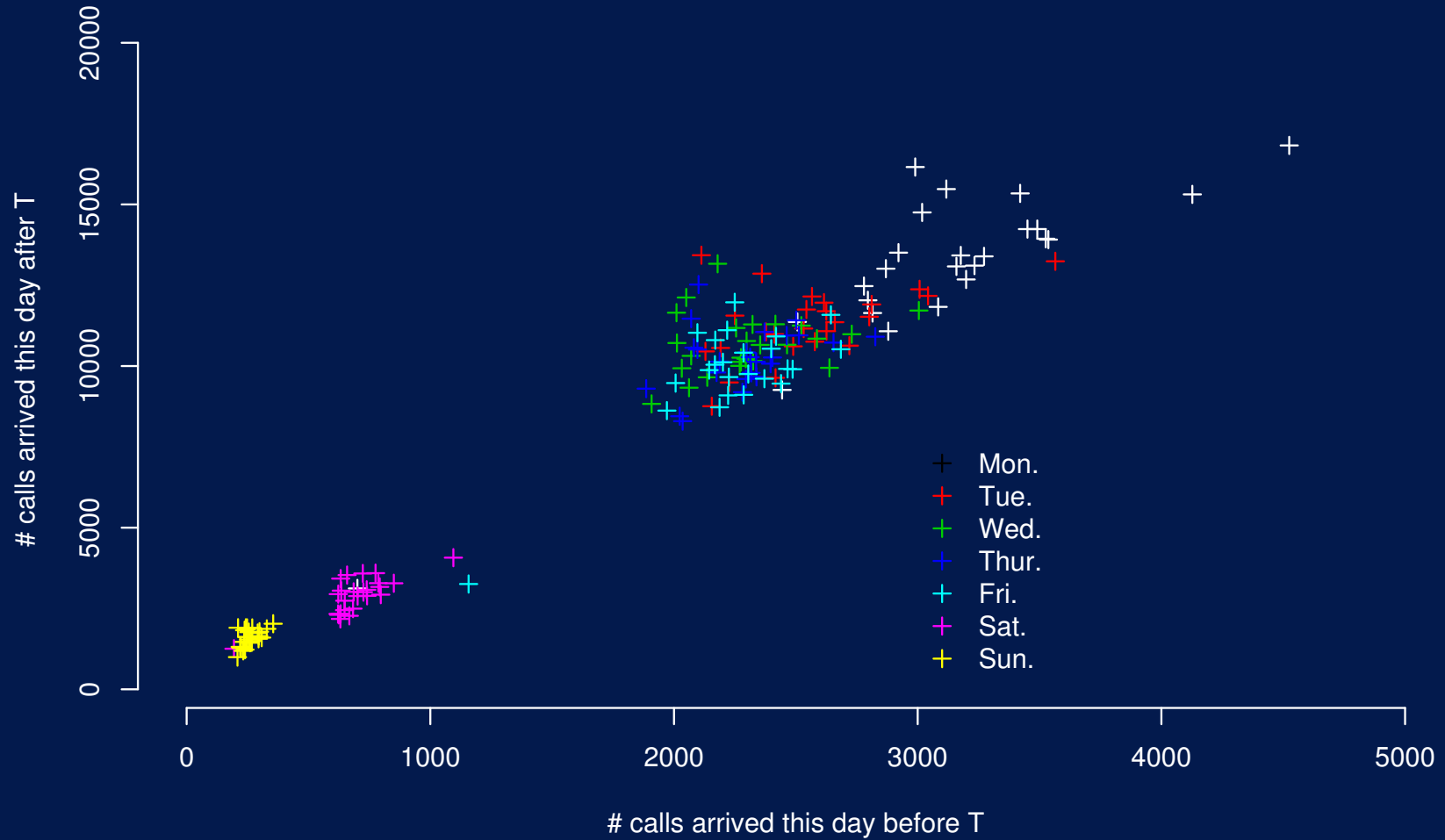
Arrival counts on different days



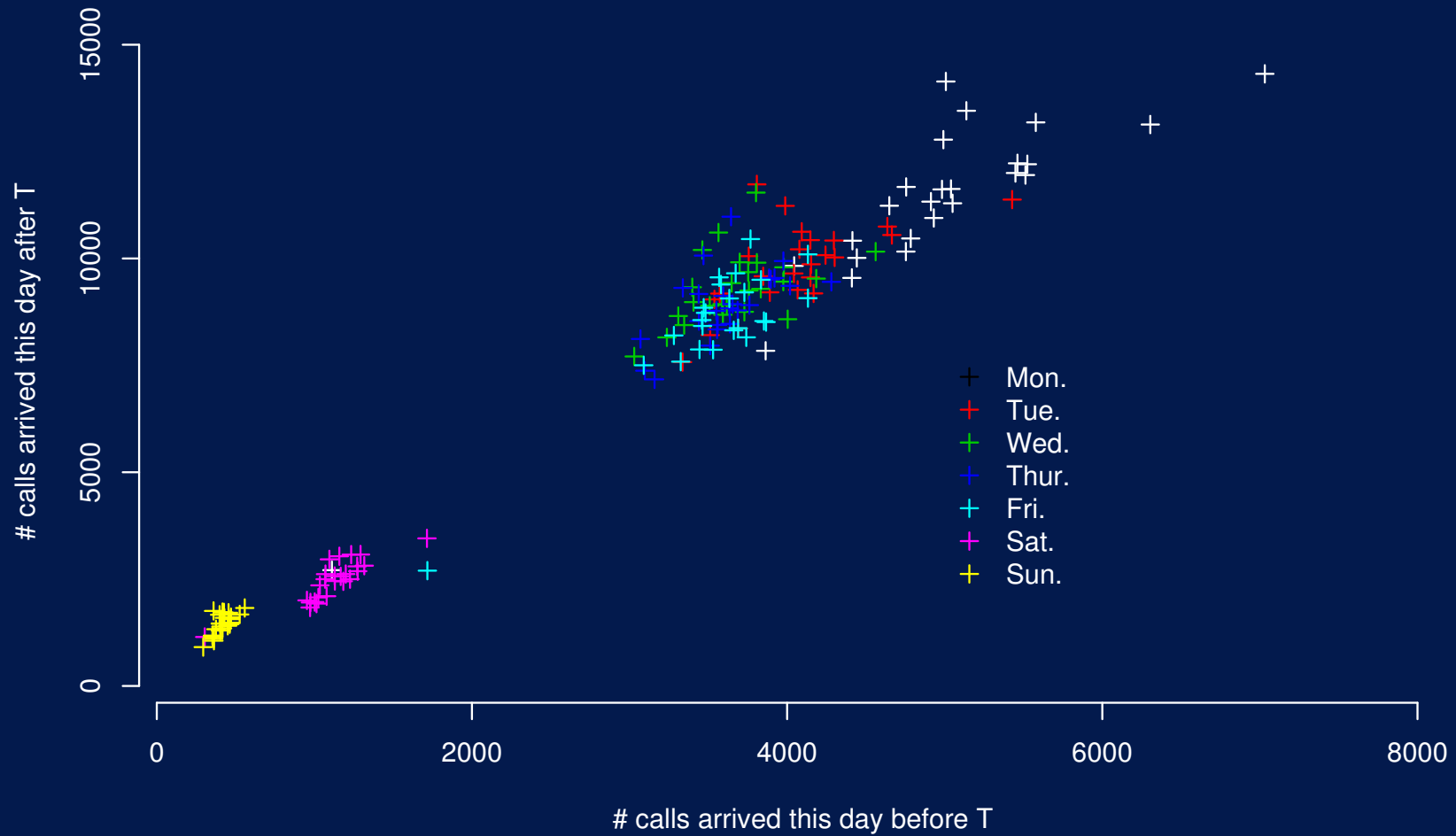
Daily residuals after removing day-effect, Box-Ljung test: $p\text{-value} = 7.534e-06$



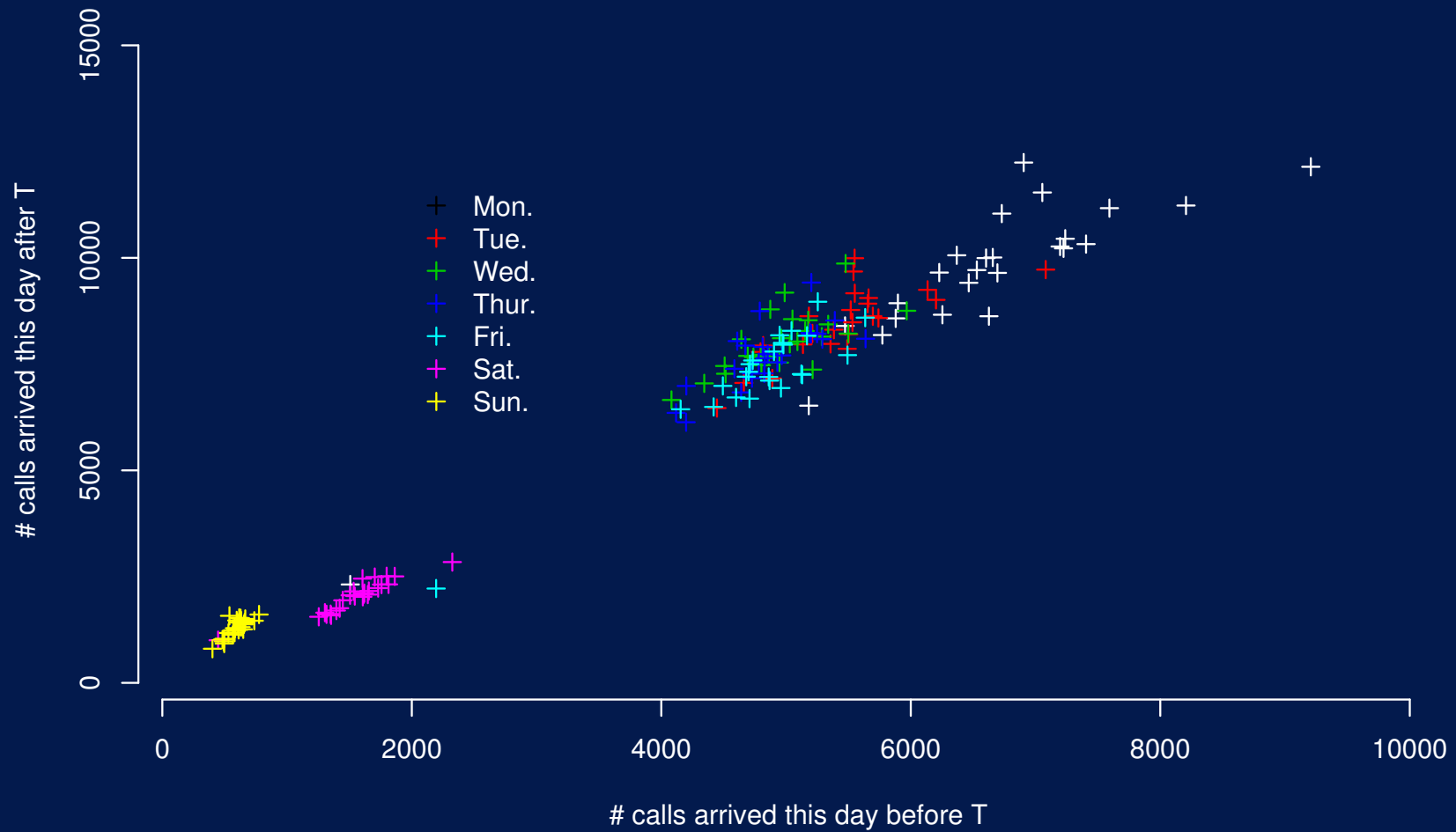
Monday, arrival counts by quarter hour



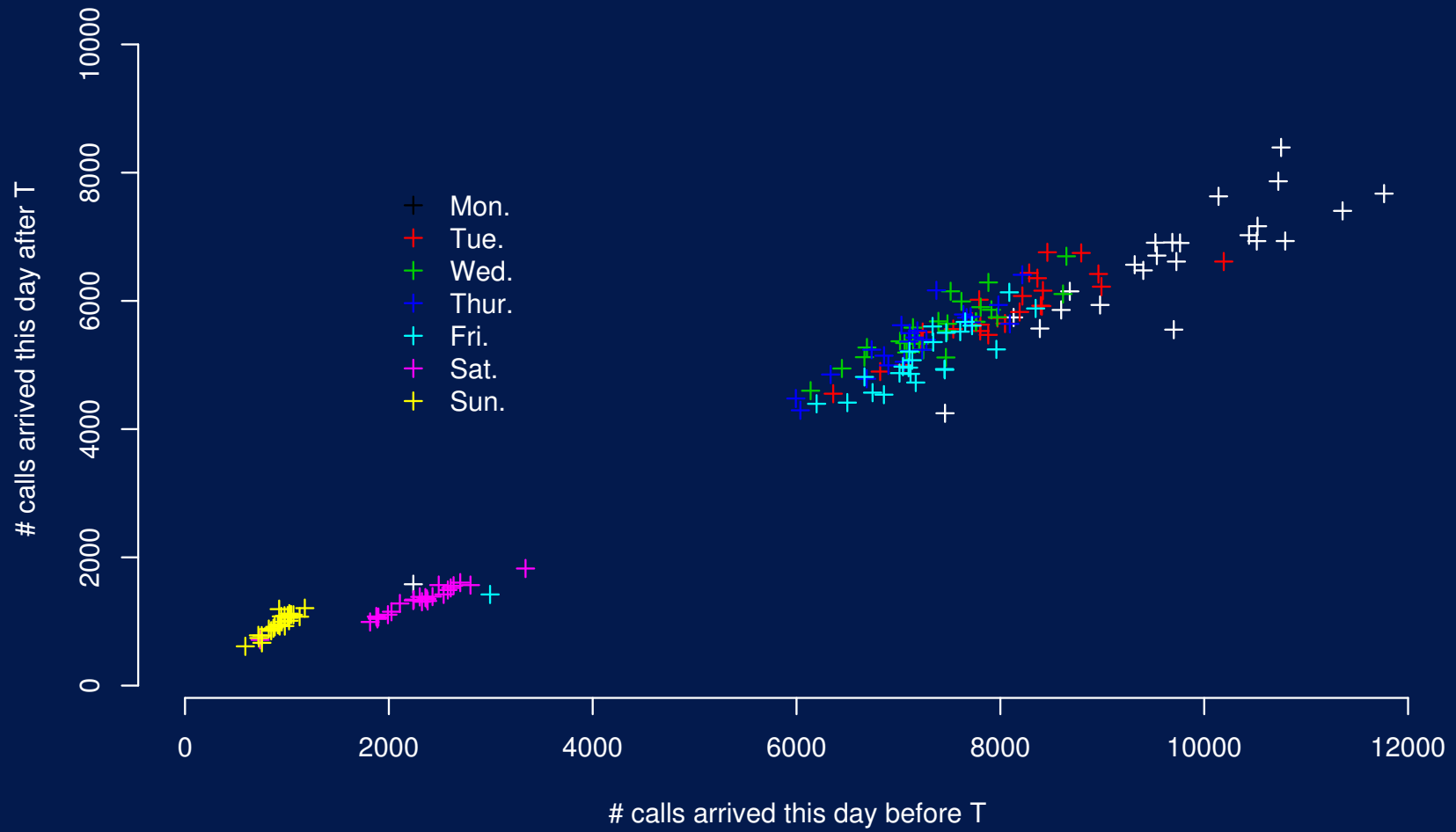
Scatter plot, all days $T=10$ am



Scatter plot, all days $T=11$ am



Scatter plot, all days $T=12$ am



Scatter plot, all days $T=2$ pm

Arrival process modeling

Arrival processes can be nonstationary, doubly stochastic, dependent, etc.

Examples, for single type: nonstationary Poisson with random inflation factor each day; arrival rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$.

(Whitt, Avramidis et al. 2004).

Arrival process modeling

Arrival processes can be nonstationary, doubly stochastic, dependent, etc.

Examples, for single type: nonstationary Poisson with random inflation factor each day; arrival rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$.

(Whitt, Avramidis et al. 2004).

Other multivariate distributions (NORTA, etc.).

Arrival process modeling

Arrival processes can be nonstationary, doubly stochastic, dependent, etc.

Examples, for single type: nonstationary Poisson with random inflation factor each day; arrival rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$.

(Whitt, Avramidis et al. 2004).

Other multivariate distributions (NORTA, etc.).

But: other relevant information can be used to forecast call volumes.

All this information should be part of our models!

Arrival process modeling

Arrival processes can be nonstationary, doubly stochastic, dependent, etc.

Examples, for single type: nonstationary Poisson with random inflation factor each day; arrival rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$.

(Whitt, Avramidis et al. 2004).

Other multivariate distributions (NORTA, etc.).

But: other relevant information can be used to forecast call volumes.

All this information should be part of our models!

Service time distributions;

Patience time distributions;

Etc.

Steady-state approximations of SL

Single skill: M/M/s (Erlang-C Analysis).

μ = service rate; λ = arrival rate; $\rho = \lambda/\mu = \text{load}$

s = staffing (number of servers)

Delay probability:

$$P(\text{wait} > 0) = \frac{\frac{\rho^s}{s!} \frac{s}{s - \rho}}{\frac{\rho^s}{s!} \frac{s}{s - \rho} + \sum_{i=0}^{s-1} \frac{\rho^i}{i!}}$$

Large system: square-root safety staffing.

Halfin-Whitt (1981) limit: sequence of M/M/s queues;

Queue n has $s_n = n$ servers, arrival rate λ_n , and load $\rho_n = \lambda_n/\mu$.

W_n = steady-state waiting time.

Large system: square-root safety staffing.

Halfin-Whitt (1981) limit: sequence of M/M/s queues;

Queue n has $s_n = n$ servers, arrival rate λ_n , and load $\rho_n = \lambda_n/\mu$.

W_n = steady-state waiting time.

Theorem. For $n \rightarrow \infty$, if $\lambda_n \rightarrow \infty$ and $(1 - \rho_n/n)\sqrt{n} \rightarrow \beta$ for $0 < \beta < 1$, then

$$P(W_n > 0) \rightarrow \delta = \frac{1}{1 + \beta\Phi(\beta)/\phi(\beta)}$$

where Φ and ϕ are the standard normal c.d.f. and p.d.f.

Large system: square-root safety staffing.

Halfin-Whitt (1981) limit: sequence of M/M/s queues;

Queue n has $s_n = n$ servers, arrival rate λ_n , and load $\rho_n = \lambda_n/\mu$.

W_n = steady-state waiting time.

Theorem. For $n \rightarrow \infty$, if $\lambda_n \rightarrow \infty$ and $(1 - \rho_n/n)\sqrt{n} \rightarrow \beta$ for $0 < \beta < 1$, then

$$P(W_n > 0) \rightarrow \delta = \frac{1}{1 + \beta\Phi(\beta)/\phi(\beta)}$$

where Φ and ϕ are the standard normal c.d.f. and p.d.f.

Square-root formula: For load ρ and delay probability δ , staff

$$n = \rho + \beta\sqrt{\rho}$$

agents.

Large system: square-root safety staffing.

Halfin-Whitt (1981) limit: sequence of M/M/s queues;

Queue n has $s_n = n$ servers, arrival rate λ_n , and load $\rho_n = \lambda_n/\mu$.

W_n = steady-state waiting time.

Theorem. For $n \rightarrow \infty$, if $\lambda_n \rightarrow \infty$ and $(1 - \rho_n/n)\sqrt{n} \rightarrow \beta$ for $0 < \beta < 1$, then

$$P(W_n > 0) \rightarrow \delta = \frac{1}{1 + \beta\Phi(\beta)/\phi(\beta)}$$

where Φ and ϕ are the standard normal c.d.f. and p.d.f.

Square-root formula: For load ρ and delay probability δ , staff

$$n = \rho + \beta\sqrt{\rho}$$

agents. Simpler than inverting M/M/s formula.

Time-dependent arrival rate.

Pointwise Stationary (PS) approximation (Green and Kolesar 1991, Whitt 1991):

1. Use steady-state approximation at each time point t ;
2. integrate result (performance) with respect to t .

Time-dependent arrival rate.

Pointwise Stationary (PS) approximation (Green and Kolesar 1991, Whitt 1991):

1. Use steady-state approximation at each time point t ;
2. integrate result (performance) with respect to t .

Improvement: use arrival rate at t for performance at $t + \delta$, for some delay δ (e.g., Ingolfsson et al.).

Multiskill case

Multiskill case

Koole and Talim (2000): static overflow routing, loss system.

Multiskill case

Koole and Talim (2000): static overflow routing, loss system.

Call type k has a Poisson arrival process, with rate λ_k , and ordered list of agent types to look for.

Multiskill case

Koole and Talim (2000): static overflow routing, loss system.

Call type k has a Poisson arrival process, with rate λ_k , and ordered list of agent types to look for.

If all these agent types are busy, the call is lost.

Multiskill case

Koole and Talim (2000): static overflow routing, loss system.

Call type k has a Poisson arrival process, with rate λ_k , and ordered list of agent types to look for.

If all these agent types are busy, the call is lost.

The net arrival process at each skill group is assumed Poisson.

Multiskill case

Koole and Talim (2000): static overflow routing, loss system.

Call type k has a Poisson arrival process, with rate λ_k , and ordered list of agent types to look for.

If all these agent types are busy, the call is lost.

The net arrival process at each skill group is assumed Poisson.

Blocking probabilities per skill group and loss rates per call type are computed by solving a system of nonlinear equations.

Other (better) approximations of loss rates:

Equivalent Random Method: two-moment approximation of overflow process (Hayward 1981, Wolff 1989, Chevalier et. al. 2004).

Other (better) approximations of loss rates:

Equivalent Random Method: two-moment approximation of overflow process (Hayward 1981, Wolff 1989, Chevalier et. al. 2004).

Hyper-exponential approximation (Franx et. al. 2005): Reportedly very successful at approximating by-type loss rates. Restrictions on routing.

Other (better) approximations of loss rates:

Equivalent Random Method: two-moment approximation of overflow process (Hayward 1981, Wolff 1989, Chevalier et. al. 2004).

Hyper-exponential approximation (Franx et. al. 2005): Reportedly very successful at approximating by-type loss rates. Restrictions on routing.

Loss-delay approximation (Avramidis et al. 2006).

Other (better) approximations of loss rates:

Equivalent Random Method: two-moment approximation of overflow process (Hayward 1981, Wolff 1989, Chevalier et. al. 2004).

Hyper-exponential approximation (Franx et. al. 2005): Reportedly very successful at approximating by-type loss rates. Restrictions on routing.

Loss-delay approximation (Avramidis et al. 2006).

Assume customer type k joins a waiting queue at the last skill group in its list, if all are busy.

Other (better) approximations of loss rates:

Equivalent Random Method: two-moment approximation of overflow process (Hayward 1981, Wolff 1989, Chevalier et. al. 2004).

Hyper-exponential approximation (Franx et. al. 2005): Reportedly very successful at approximating by-type loss rates. Restrictions on routing.

Loss-delay approximation (Avramidis et al. 2006).

Assume customer type k joins a waiting queue at the last skill group in its list, if all are busy.

The SL for type k is approximated using a birth-and-death CTMC model at that queue, for each k .

Other (better) approximations of loss rates:

Equivalent Random Method: two-moment approximation of overflow process (Hayward 1981, Wolff 1989, Chevalier et. al. 2004).

Hyper-exponential approximation (Franx et. al. 2005): Reportedly very successful at approximating by-type loss rates. Restrictions on routing.

Loss-delay approximation (Avramidis et al. 2006).

Assume customer type k joins a waiting queue at the last skill group in its list, if all are busy.

The SL for type k is approximated using a birth-and-death CTMC model at that queue, for each k .

Allows abandonments and general routing.

Other (better) approximations of loss rates:

Equivalent Random Method: two-moment approximation of overflow process (Hayward 1981, Wolff 1989, Chevalier et. al. 2004).

Hyper-exponential approximation (Franx et. al. 2005): Reportedly very successful at approximating by-type loss rates. Restrictions on routing.

Loss-delay approximation (Avramidis et al. 2006).

Assume customer type k joins a waiting queue at the last skill group in its list, if all are busy.

The SL for type k is approximated using a birth-and-death CTMC model at that queue, for each k .

Allows abandonments and general routing.

Not very realistic, but useful for rough-cut first-step in staffing optimization.

Optimization: Integer Programming Formulation

Call types	$k = 1, \dots, K;$
Agent types (or skill groups)	$i = 1, \dots, I;$
Periods	$p = 1, \dots, P;$
Shift types	$q = 1, \dots, Q.$

Optimization: Integer Programming Formulation

Call types $k = 1, \dots, K;$

Agent types (or skill groups) $i = 1, \dots, I;$

Periods $p = 1, \dots, P;$

Shift types $q = 1, \dots, Q.$

The shift specifies the time when the agent starts working, when he/she finishes, and all the lunch and coffee breaks.

Optimization: Integer Programming Formulation

Call types $k = 1, \dots, K;$

Agent types (or skill groups) $i = 1, \dots, I;$

Periods $p = 1, \dots, P;$

Shift types $q = 1, \dots, Q.$

The shift specifies the time when the agent starts working, when he/she finishes, and all the lunch and coffee breaks.

Costs: $\mathbf{c} = (c_{1,1}, \dots, c_{1,Q}, \dots, c_{I,1}, \dots, c_{I,Q})$ where $c_{i,q}$ = cost of an agent of type i having shift q .

Optimization: Integer Programming Formulation

Call types	$k = 1, \dots, K;$
Agent types (or skill groups)	$i = 1, \dots, I;$
Periods	$p = 1, \dots, P;$
Shift types	$q = 1, \dots, Q.$

The shift specifies the time when the agent starts working, when he/she finishes, and all the lunch and coffee breaks.

Costs: $\mathbf{c} = (c_{1,1}, \dots, c_{1,Q}, \dots, c_{I,1}, \dots, c_{I,Q})$ where
 $c_{i,q}$ = cost of an agent of type i having shift q .

Decision variables: $\mathbf{x} = (x_{1,1}, \dots, x_{1,Q}, \dots, x_{I,1}, \dots, x_{I,Q})$ where
 $x_{i,q}$ = number of agents of type i having shift q .

Optimization: Integer Programming Formulation

Call types	$k = 1, \dots, K;$
Agent types (or skill groups)	$i = 1, \dots, I;$
Periods	$p = 1, \dots, P;$
Shift types	$q = 1, \dots, Q.$

The shift specifies the time when the agent starts working, when he/she finishes, and all the lunch and coffee breaks.

Costs: $\mathbf{c} = (c_{1,1}, \dots, c_{1,Q}, \dots, c_{I,1}, \dots, c_{I,Q})$ where
 $c_{i,q}$ = cost of an agent of type i having shift q .

Decision variables: $\mathbf{x} = (x_{1,1}, \dots, x_{1,Q}, \dots, x_{I,1}, \dots, x_{I,Q})$ where
 $x_{i,q}$ = number of agents of type i having shift q .

Auxiliary variables: $\mathbf{y} = (y_{1,1}, \dots, y_{1,P}, \dots, y_{I,1}, \dots, y_{I,P})$ where
 $y_{i,p}$ = number of agents of type i in period p .

Optimization: Integer Programming Formulation

Call types	$k = 1, \dots, K;$
Agent types (or skill groups)	$i = 1, \dots, I;$
Periods	$p = 1, \dots, P;$
Shift types	$q = 1, \dots, Q.$

The shift specifies the time when the agent starts working, when he/she finishes, and all the lunch and coffee breaks.

Costs: $\mathbf{c} = (c_{1,1}, \dots, c_{1,Q}, \dots, c_{I,1}, \dots, c_{I,Q})$ where
 $c_{i,q}$ = cost of an agent of type i having shift q .

Decision variables: $\mathbf{x} = (x_{1,1}, \dots, x_{1,Q}, \dots, x_{I,1}, \dots, x_{I,Q})$ where
 $x_{i,q}$ = number of agents of type i having shift q .

Auxiliary variables: $\mathbf{y} = (y_{1,1}, \dots, y_{1,P}, \dots, y_{I,1}, \dots, y_{I,P})$ where
 $y_{i,p}$ = number of agents of type i in period p .

We have $\mathbf{y} = \mathbf{A}\mathbf{x}$ where \mathbf{A} is block diagonal with i blocks $\tilde{\mathbf{A}}$, and element (p, q) of $\tilde{\mathbf{A}}$ is 1 if shift q covers period p , 0 otherwise.

$g_{k,p}(\mathbf{y})$ = long-run SL for call type k in period p .

E.g., fraction of calls answered within 20 seconds, or $s_{k,p}$ seconds:

$$g_{k,p}(\mathbf{y}) = \frac{E[\text{num. of calls in period } p \text{ answered within the limit}]}{E[\text{num. of arrivals in period } p]}.$$

$g_{k,p}(\mathbf{y})$ = long-run SL for call type k in period p .

E.g., fraction of calls answered within 20 seconds, or $s_{k,p}$ seconds:

$$g_{k,p}(\mathbf{y}) = \frac{E[\text{num. of calls in period } p \text{ answered within the limit}]}{E[\text{num. of arrivals in period } p]}.$$

$g_p(\mathbf{y})$ = aggregated long-run SL over period p .

$g_k(\mathbf{y})$ = aggregated long-run SL for call type k .

$g(\mathbf{y})$ = aggregated long-run SL overall.

$g_{k,p}(\mathbf{y})$ = long-run SL for call type k in period p .

E.g., fraction of calls answered within 20 seconds, or $s_{k,p}$ seconds:

$$g_{k,p}(\mathbf{y}) = \frac{E[\text{num. of calls in period } p \text{ answered within the limit}]}{E[\text{num. of arrivals in period } p]}.$$

$g_p(\mathbf{y})$ = aggregated long-run SL over period p .

$g_k(\mathbf{y})$ = aggregated long-run SL for call type k .

$g(\mathbf{y})$ = aggregated long-run SL overall.

These functions g_\bullet (ratios of expectations) are unknown and often very complicated. Each may depend on entire vector \mathbf{y} . They can be either

- approximated via simplified queueing models, or
- estimated by simulation.

Here, the routing rules are assumed **fixed**; we do not optimize them.

Scheduling problem

$$\min \quad \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q}$$

subject to

$$\begin{aligned} \mathbf{A} \mathbf{x} &= \mathbf{y}, \\ g_{k,p}(\mathbf{y}) &\geq l_{k,p} && \text{for all } k, p, \\ g_p(\mathbf{y}) &\geq l_p && \text{for all } p, \\ g_k(\mathbf{y}) &\geq l_k && \text{for all } k, \\ g(\mathbf{y}) &\geq l, \\ \mathbf{x} &\geq 0, \text{ and integer.} \end{aligned}$$

Staffing Problem

Relaxation: forget about the admissibility of schedules; just assume that any staffing is admissible.

Staffing Problem

Relaxation: forget about the admissibility of schedules; just assume that any staffing is admissible.

Costs: $\mathbf{c} = (c_{1,1}, \dots, c_{1,P}, \dots, c_{I,1}, \dots, c_{I,P})$ where
 $c_{i,p}$ = cost of an agent of type i in period p .

Staffing Problem

Relaxation: forget about the admissibility of schedules; just assume that any staffing is admissible.

Costs: $\mathbf{c} = (c_{1,1}, \dots, c_{1,P}, \dots, c_{I,1}, \dots, c_{I,P})$ where $c_{i,p}$ = cost of an agent of type i in period p .

$$\begin{aligned} \min \quad & \mathbf{c}^t \mathbf{y} = \sum_{i=1}^I \sum_{p=1}^P c_{i,p} y_{i,p} \\ \text{subject to} \quad & g_{k,p}(\mathbf{y}) \geq l_{k,p} \quad \text{for all } k, p, \\ & g_p(\mathbf{y}) \geq l_p \quad \text{for all } p, \\ & g_k(\mathbf{y}) \geq l_k \quad \text{for all } k, \\ & g(\mathbf{y}) \geq l, \\ & \mathbf{y} \geq 0, \text{ and integer.} \end{aligned}$$

Single period, steady-state approximation

Take one period at a time and assume that the system is in steady-state.

Single period, steady-state approximation

Take one period at a time and assume that the system is in steady-state.

$\mathbf{c} = (c_1, \dots, c_I)$ where c_i = cost of an agent of type i .

Single period, steady-state approximation

Take one period at a time and assume that the system is in steady-state.

$\mathbf{c} = (c_1, \dots, c_I)$ where c_i = cost of an agent of type i .

$\mathbf{y} = (y_1, \dots, y_I)$ where y_i = number of agents of type i .

Single period, steady-state approximation

Take one period at a time and assume that the system is in steady-state.

$\mathbf{c} = (c_1, \dots, c_I)$ where c_i = cost of an agent of type i .

$\mathbf{y} = (y_1, \dots, y_I)$ where y_i = number of agents of type i .

$$\begin{aligned} \min \quad & \mathbf{c}^t \mathbf{y} = \sum_{i=1}^I c_i y_i \\ \text{subject to} \quad & g_k(\mathbf{y}) \geq l_k \quad \text{for all } k, \\ & g(\mathbf{y}) \geq l, \\ & \mathbf{y} \geq 0, \text{ and integer.} \end{aligned}$$

Single period, steady-state approximation

Take one period at a time and assume that the system is in steady-state.

$\mathbf{c} = (c_1, \dots, c_I)$ where c_i = cost of an agent of type i .

$\mathbf{y} = (y_1, \dots, y_I)$ where y_i = number of agents of type i .

$$\begin{aligned} \min \quad & \mathbf{c}^t \mathbf{y} = \sum_{i=1}^I c_i y_i \\ \text{subject to} \quad & g_k(\mathbf{y}) \geq l_k \quad \text{for all } k, \\ & g(\mathbf{y}) \geq l, \\ & \mathbf{y} \geq 0, \text{ and integer.} \end{aligned}$$

Neglects the transient and overflow effect across periods.

Scheduling by a two-step method

First step: solve the staffing problem (easier).

Let y^* be an optimal solution.

Scheduling by a two-step method

First step: solve the staffing problem (easier).

Let \mathbf{y}^* be an optimal solution.

Second step: find an admissible schedule that covers the staffing \mathbf{y}^* , by solving:

$$\begin{aligned} \min \quad & \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q} \\ \text{subject to} \quad & \mathbf{A} \mathbf{x} \geq \mathbf{y}^*, \\ & \mathbf{x} \geq 0, \text{ and integer.} \end{aligned}$$

Scheduling by a two-step method

First step: solve the staffing problem (easier).

Let \mathbf{y}^* be an optimal solution.

Second step: find an admissible schedule that covers the staffing \mathbf{y}^* , by solving:

$$\begin{aligned} \min \quad & \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q} \\ \text{subject to} \quad & \mathbf{A} \mathbf{x} \geq \mathbf{y}^*, \\ & \mathbf{x} \geq 0, \text{ and integer.} \end{aligned}$$

This does not give an optimal solution to the original scheduling problem.
And the gap could be significant.

Scheduling by a two-step method

First step: solve the staffing problem (easier).

Let \mathbf{y}^* be an optimal solution.

Second step: find an admissible schedule that covers the staffing \mathbf{y}^* , by solving:

$$\begin{aligned} \min \quad & \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q} \\ \text{subject to} \quad & \mathbf{A} \mathbf{x} \geq \mathbf{y}^*, \\ & \mathbf{x} \geq 0, \text{ and integer.} \end{aligned}$$

This does not give an optimal solution to the original scheduling problem.

And the gap could be significant.

Potential discrepancy in mix of agent types across successive periods.

Scheduling by a two-step method

First step: solve the staffing problem (easier).

Let \mathbf{y}^* be an optimal solution.

Second step: find an admissible schedule that covers the staffing \mathbf{y}^* , by solving:

$$\begin{aligned} \min \quad & \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q} \\ \text{subject to} \quad & \mathbf{A} \mathbf{x} \geq \mathbf{y}^*, \\ & \mathbf{x} \geq 0, \text{ and integer.} \end{aligned}$$

This does not give an optimal solution to the original scheduling problem.

And the gap could be significant.

Potential discrepancy in mix of agent types across successive periods.

Can be alleviated by skill-transfer decision variables.

Skill-transfer decision variables (Bhulai, Koole, and Pot 2005):

$z_{p,j,i}$ = number of agents of type j that work as type- i agents in period p (they use only part of their skills).

$$\min \quad \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q}$$

subject to

$$\sum_{q: \tilde{A}_{p,q}=1} x_{i,q} + \sum_{j \in \mathcal{S}_i^+} z_{p,j,i} - \sum_{j \in \mathcal{S}_i^-} z_{p,i,j} \geq y_{i,p} \quad \forall p, i$$

all $x_{i,q}, z_{p,j,i} \geq 0$ and integer.

Sample-path optimization via simulation

[Atlason, Epelman, and Henderson, 2004; Cezik and L'Ecuyer 2005]

We simulate n independent operating days (could also be weeks, etc.) of the center, to estimate the functions g_{\bullet} .

Sample-path optimization via simulation

[Atlason, Epelman, and Henderson, 2004; Cezik and L'Ecuyer 2005]

We simulate n independent operating days (could also be weeks, etc.) of the center, to estimate the functions g_{\bullet} .

Let ω represent the source of randomness, i.e., the sequence of independent uniform r.v.'s underlying the entire simulation (n runs).

Sample-path optimization via simulation

[Atlason, Epelman, and Henderson, 2004; Cezik and L'Ecuyer 2005]

We simulate n independent operating days (could also be weeks, etc.) of the center, to estimate the functions g_{\bullet} .

Let ω represent the source of randomness, i.e., the sequence of independent uniform r.v.'s underlying the entire simulation (n runs).

For a ω , the empirical SL's over the n simulation runs are:

$G_{n,k,p}(\mathbf{y}, \omega)$ for call type k in period p ;

$G_{n,p}(\mathbf{y}, \omega)$ aggregated over period p ;

$G_{n,k}(\mathbf{y}, \omega)$ aggregated for call type k ;

$G_n(\mathbf{y}, \omega)$ aggregated overall.

Sample-path optimization via simulation

[Atlason, Epelman, and Henderson, 2004; Cezik and L'Ecuyer 2005]

We simulate n independent operating days (could also be weeks, etc.) of the center, to estimate the functions g_{\bullet} .

Let ω represent the source of randomness, i.e., the sequence of independent uniform r.v.'s underlying the entire simulation (n runs).

For a ω , the empirical SL's over the n simulation runs are:

$G_{n,k,p}(\mathbf{y}, \omega)$ for call type k in period p ;

$G_{n,p}(\mathbf{y}, \omega)$ aggregated over period p ;

$G_{n,k}(\mathbf{y}, \omega)$ aggregated for call type k ;

$G_n(\mathbf{y}, \omega)$ aggregated overall.

For a fixed ω , these are deterministic functions of \mathbf{y} .

Sample-path optimization via simulation

[Atlason, Epelman, and Henderson, 2004; Cezik and L'Ecuyer 2005]

We simulate n independent operating days (could also be weeks, etc.) of the center, to estimate the functions g_{\bullet} .

Let ω represent the source of randomness, i.e., the sequence of independent uniform r.v.'s underlying the entire simulation (n runs).

For a ω , the empirical SL's over the n simulation runs are:

$G_{n,k,p}(\mathbf{y}, \omega)$ for call type k in period p ;

$G_{n,p}(\mathbf{y}, \omega)$ aggregated over period p ;

$G_{n,k}(\mathbf{y}, \omega)$ aggregated for call type k ;

$G_n(\mathbf{y}, \omega)$ aggregated overall.

For a fixed ω , these are deterministic functions of \mathbf{y} .

We replace the functions g by the estimators G and optimize.

Sample-path optimization via simulation

[Atlason, Epelman, and Henderson, 2004; Cezik and L'Ecuyer 2005]

We simulate n independent operating days (could also be weeks, etc.) of the center, to estimate the functions g_{\bullet} .

Let ω represent the source of randomness, i.e., the sequence of independent uniform r.v.'s underlying the entire simulation (n runs).

For a ω , the empirical SL's over the n simulation runs are:

$G_{n,k,p}(\mathbf{y}, \omega)$ for call type k in period p ;

$G_{n,p}(\mathbf{y}, \omega)$ aggregated over period p ;

$G_{n,k}(\mathbf{y}, \omega)$ aggregated for call type k ;

$G_n(\mathbf{y}, \omega)$ aggregated overall.

For a fixed ω , these are deterministic functions of \mathbf{y} .

We replace the functions g by the estimators G and optimize.

To compute them at different values of \mathbf{y} , we simply use simulation with common random numbers.

Empirical scheduling optimization problem (sample version of the problem):

$$\begin{aligned} \min \quad & \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q} \\ \text{subject to} \quad & \mathbf{Ax} = \mathbf{y}, \\ & G_{n,k,p}(\mathbf{y}) \geq l_{k,p} \quad \text{for all } k, p, \\ & G_{n,p}(\mathbf{y}) \geq l_p \quad \text{for all } p, \\ & G_{n,k}(\mathbf{y}) \geq l_k \quad \text{for all } k, \\ & G_n(\mathbf{y}) \geq l, \\ & \mathbf{x} \geq 0, \text{ and integer.} \end{aligned}$$

Empirical scheduling optimization problem (sample version of the problem):

$$\begin{aligned}
 \min \quad & \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q} \\
 \text{subject to} \quad & \mathbf{A} \mathbf{x} = \mathbf{y}, \\
 & G_{n,k,p}(\mathbf{y}) \geq l_{k,p} \quad \text{for all } k, p, \\
 & G_{n,p}(\mathbf{y}) \geq l_p \quad \text{for all } p, \\
 & G_{n,k}(\mathbf{y}) \geq l_k \quad \text{for all } k, \\
 & G_n(\mathbf{y}) \geq l, \\
 & \mathbf{x} \geq 0, \text{ and integer.}
 \end{aligned}$$

Under mild conditions, when $n \rightarrow \infty$, the optimal solution of the sample problem converges w.p.1 to that of the original problem.

Empirical scheduling optimization problem (sample version of the problem):

$$\begin{aligned}
 \min \quad & \mathbf{c}^t \mathbf{x} = \sum_{i=1}^I \sum_{q=1}^Q c_{i,q} x_{i,q} \\
 \text{subject to} \quad & \mathbf{A} \mathbf{x} = \mathbf{y}, \\
 & G_{n,k,p}(\mathbf{y}) \geq l_{k,p} \quad \text{for all } k, p, \\
 & G_{n,p}(\mathbf{y}) \geq l_p \quad \text{for all } p, \\
 & G_{n,k}(\mathbf{y}) \geq l_k \quad \text{for all } k, \\
 & G_n(\mathbf{y}) \geq l, \\
 & \mathbf{x} \geq 0, \text{ and integer.}
 \end{aligned}$$

Under mild conditions, when $n \rightarrow \infty$, the optimal solution of the sample problem converges w.p.1 to that of the original problem.

Similar formulation for the [staffing](#) problem.

Neighborhood search for staffing problem

[Avramidis et al. 2006]

Stage 1: Start with feasible solution and decrease cost using neighborhood search with “loss-delay” approximation to evaluate SL and select next solution.

Neighborhood search for staffing problem

[Avramidis et al. 2006]

Stage 1: Start with feasible solution and decrease cost using neighborhood search with “loss-delay” approximation to evaluate SL and select next solution.

Outputs a locally optimal solution (w.r.t. loss-delay approx.).

Neighborhood search for staffing problem

[Avramidis et al. 2006]

Stage 1: Start with feasible solution and decrease cost using neighborhood search with “loss-delay” approximation to evaluate SL and select next solution.

Outputs a locally optimal solution (w.r.t. loss-delay approx.).

Stage 2: simulation-based local adjustment to stage-1 solution to account for two mutually exclusive cases:

- the incumbent is infeasible and we seek a nearby feasible solution;
- the incumbent is feasible and we seek further cost reduction.

Variance reduction: examples

Suppose the arrival process is Poisson with rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$ and $\lambda(\cdot)$ is deterministic.

Variance reduction: examples

Suppose the arrival process is Poisson with rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$ and $\lambda(\cdot)$ is deterministic.

$X_i = G_i(s_0)$ = number of calls who waited less than s_0 seconds on day i . Crude estimator of $\mu = \mathbb{E}[X_i]$:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

Variance reduction: examples

Suppose the arrival process is Poisson with rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$ and $\lambda(\cdot)$ is deterministic.

$X_i = G_i(s_0)$ = number of calls who waited less than s_0 seconds on day i . Crude estimator of $\mu = \mathbb{E}[X_i]$:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

The realization B_i of B has a large impact on $G(s_0)$, so we may want to use quasi-Monte Carlo or stratification w.r.t. B .

Variance reduction: examples

Suppose the arrival process is Poisson with rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$ and $\lambda(\cdot)$ is deterministic.

$X_i = G_i(s_0)$ = number of calls who waited less than s_0 seconds on day i . Crude estimator of $\mu = \mathbb{E}[X_i]$:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

The realization B_i of B has a large impact on $G(s_0)$, so we may want to use quasi-Monte Carlo or stratification w.r.t. B .

Suppose $B = F_B^{-1}(U)$ where $U \sim \text{Unif}(0, 1)$.

Variance reduction: examples

Suppose the arrival process is Poisson with rate $B\lambda(t)$ at time t , where $\mathbb{E}[B] = 1$ and $\lambda(\cdot)$ is deterministic.

$X_i = G_i(s_0)$ = number of calls who waited less than s_0 seconds on day i . Crude estimator of $\mu = \mathbb{E}[X_i]$:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

The realization B_i of B has a large impact on $G(s_0)$, so we may want to use quasi-Monte Carlo or stratification w.r.t. B .

Suppose $B = F_B^{-1}(U)$ where $U \sim \text{Unif}(0, 1)$.

Stratification on U in k strata: for $s = 1, \dots, k$, generate n_s indep. r.v.'s $U^{(s)} \sim \text{Unif}((s-1)/k, s/k)$ and denote $X_{s,1}, \dots, X_{s,n_s}$ the corresponding observations of $G(s_0)$.

Stratified estimator:

$$\bar{X}_{s,n} = \frac{1}{k} \sum_{s=1}^k \frac{1}{n_s} \sum_{i=1}^{n_s} X_{s,i}.$$

Stratified estimator:

$$\bar{X}_{s,n} = \frac{1}{k} \sum_{s=1}^k \frac{1}{n_s} \sum_{i=1}^{n_s} X_{s,i}.$$

The variance

$$\text{Var}[\bar{X}_{s,n}] = \frac{1}{k^2} \sum_{s=1}^k \sigma_s^2 / n_s$$

where $\sigma_s^2 = \text{Var}[X | S = s]$ (must be estimated), is minimized by taking (if we neglect rounding)

$$n_s = n_s^* = n \frac{\sigma_s}{\sum_{l=1}^k \sigma_l / k}.$$

Stratified estimator:

$$\bar{X}_{s,n} = \frac{1}{k} \sum_{s=1}^k \frac{1}{n_s} \sum_{i=1}^{n_s} X_{s,i}.$$

The variance

$$\text{Var}[\bar{X}_{s,n}] = \frac{1}{k^2} \sum_{s=1}^k \sigma_s^2 / n_s$$

where $\sigma_s^2 = \text{Var}[X | S = s]$ (must be estimated), is minimized by taking (if we neglect rounding)

$$n_s = n_s^* = n \frac{\sigma_s}{\sum_{l=1}^k \sigma_l / k}.$$

We have

$$\text{Var}[\bar{X}_n] = \text{Var}[\bar{X}_{so,n}] + \frac{1}{nk} \sum_{s=1}^k (\sigma_s - \bar{\sigma})^2 + \frac{1}{nk} \sum_{s=1}^k (\mu_s - \mu)^2.$$

Combining with a control variate.

Let A = number of arrivals during the day. CV estimator conditional on $U = u$:

$$X_c(u) = X - \beta(u)[A - \mathbb{E}[A \mid U = u]] = X - \beta(u)C(u).$$

Combining with a control variate.

Let A = number of arrivals during the day. CV estimator conditional on $U = u$:

$$X_c(u) = X - \beta(u)[A - \mathbb{E}[A | U = u]] = X - \beta(u)C(u).$$

Optimal CV coefficient depend on u :

$$\beta^*(u) = \frac{\text{Cov}[A, X | U = u]}{\text{Var}[A | U = u]} = \frac{E[C \cdot X | U = u]}{E[C^2 | U = u]}.$$

Combining with a control variate.

Let A = number of arrivals during the day. CV estimator conditional on $U = u$:

$$X_c(u) = X - \beta(u)[A - \mathbb{E}[A | U = u]] = X - \beta(u)C(u).$$

Optimal CV coefficient depend on u :

$$\beta^*(u) = \frac{\text{Cov}[A, X | U = u]}{\text{Var}[A | U = u]} = \frac{E[C \cdot X | U = u]}{E[C^2 | U = u]}.$$

Can approximate β^* by approximating the two functions

$$q_1(u) = E[C \cdot X | U = u] \text{ and } q_2(u) = E[C^2 | U = u].$$

Combining with a control variate.

Let A = number of arrivals during the day. CV estimator conditional on $U = u$:

$$X_c(u) = X - \beta(u)[A - \mathbb{E}[A | U = u]] = X - \beta(u)C(u).$$

Optimal CV coefficient depend on u :

$$\beta^*(u) = \frac{\text{Cov}[A, X | U = u]}{\text{Var}[A | U = u]} = \frac{E[C \cdot X | U = u]}{E[C^2 | U = u]}.$$

Can approximate β^* by approximating the two functions

$$q_1(u) = E[C \cdot X | U = u] \text{ and } q_2(u) = E[C^2 | U = u].$$

Use pilot runs and least squares polynomial or spline approximation.

Combining with a control variate.

Let A = number of arrivals during the day. CV estimator conditional on $U = u$:

$$X_c(u) = X - \beta(u)[A - \mathbb{E}[A | U = u]] = X - \beta(u)C(u).$$

Optimal CV coefficient depend on u :

$$\beta^*(u) = \frac{\text{Cov}[A, X | U = u]}{\text{Var}[A | U = u]} = \frac{E[C \cdot X | U = u]}{E[C^2 | U = u]}.$$

Can approximate β^* by approximating the two functions

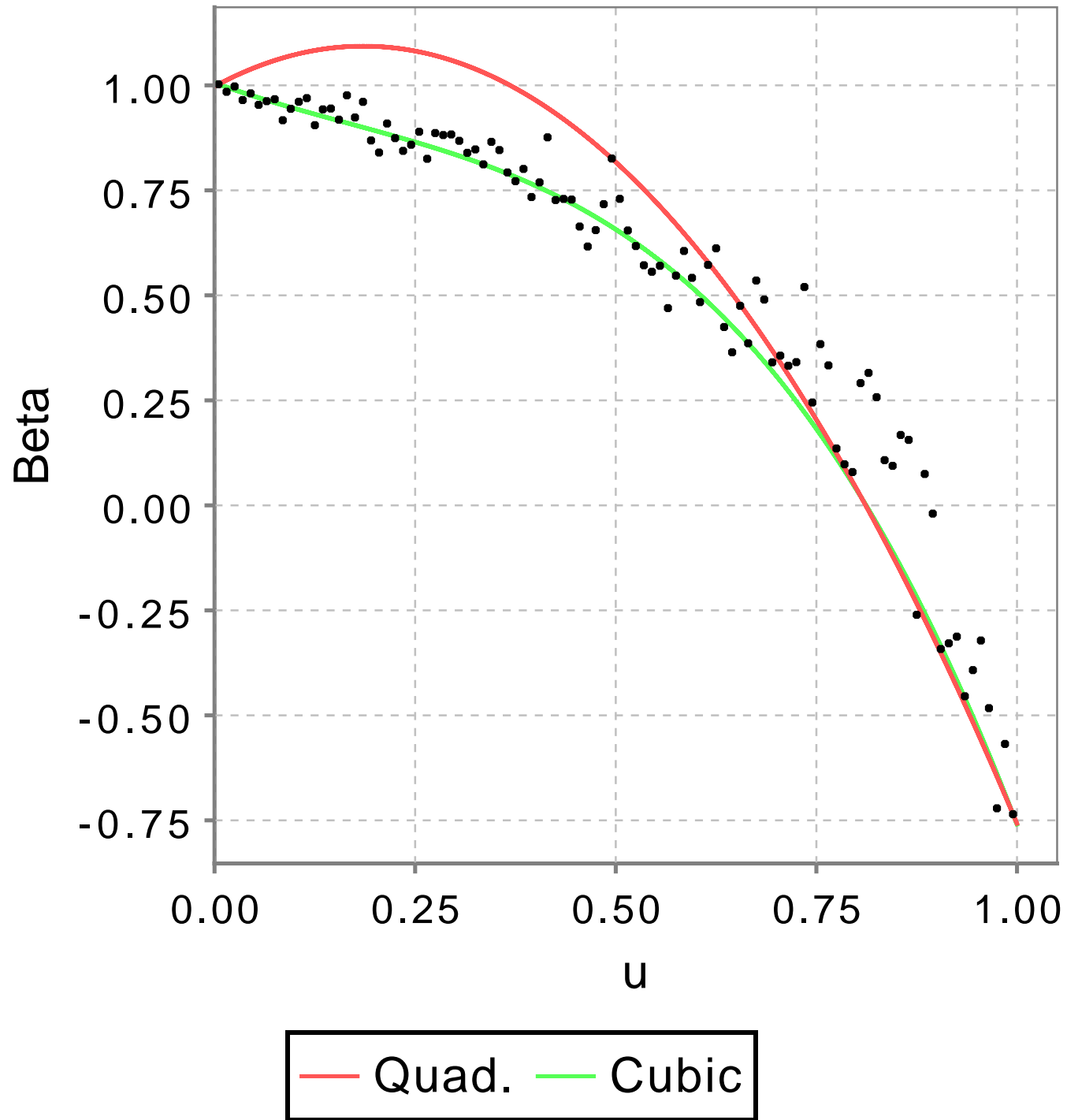
$$q_1(u) = E[C \cdot X | U = u] \text{ and } q_2(u) = E[C^2 | U = u].$$

Use pilot runs and least squares polynomial or spline approximation.

Variance in stratum s is now

$$\sigma_s^2 = \text{Var}[X_c(U^{(s)})] = \frac{1}{k} \int_{(s-1)/k}^{s/k} \text{Var}[X - \beta^*(u)C | U = u] du$$

where $U^{(s)} \sim \text{Unif}((s-1)/k, s/k)$. So the optimal allocation changes!



The function $Q^*(u)$ is fitted to the data points by a quadratic (red) and a cubic (green) curve.

Another example: Sensitivity w.r.t. mean service time parameter.

Another example: Sensitivity w.r.t. mean service time parameter.

Let $X_{1,i}$ and $X_{2,i}$ be the realizations of $G(s_0)$ with mean service time θ and $\theta + \delta$.

Another example: Sensitivity w.r.t. mean service time parameter.

Let $X_{1,i}$ and $X_{2,i}$ be the realizations of $G(s_0)$ with mean service time θ and $\theta + \delta$.

Finite difference estimator:

$$\Delta_i = [X_{2,i} - X_{1,i}] / \delta.$$

Another example: Sensitivity w.r.t. mean service time parameter.

Let $X_{1,i}$ and $X_{2,i}$ be the realizations of $G(s_0)$ with mean service time θ and $\theta + \delta$.

Finite difference estimator:

$$\Delta_i = [X_{2,i} - X_{1,i}]/\delta.$$

With IRN: $\text{Var}[\Delta_i] = O(\delta^{-2})$.

Another example: Sensitivity w.r.t. mean service time parameter.

Let $X_{1,i}$ and $X_{2,i}$ be the realizations of $G(s_0)$ with mean service time θ and $\theta + \delta$.

Finite difference estimator:

$$\Delta_i = [X_{2,i} - X_{1,i}] / \delta.$$

With IRN: $\text{Var}[\Delta_i] = O(\delta^{-2})$.

With CRN, if $G(s_0)$ was continuous, we would have $\text{Var}[\Delta_i] = O(1)$.

Another example: Sensitivity w.r.t. mean service time parameter.

Let $X_{1,i}$ and $X_{2,i}$ be the realizations of $G(s_0)$ with mean service time θ and $\theta + \delta$.

Finite difference estimator:

$$\Delta_i = [X_{2,i} - X_{1,i}] / \delta.$$

With IRN: $\text{Var}[\Delta_i] = O(\delta^{-2})$.

With CRN, if $G(s_0)$ was continuous, we would have $\text{Var}[\Delta_i] = O(1)$.

With well-synchronized CRN, we can prove that $\text{Var}[\Delta_i] = O(\delta^{-1-\epsilon})$.

Ongoing and Future Work in Our Group

- Develop a Java library for [simulation](#) of call centers.
- Variance reduction methods for simulation.
- Integrate with [optimization](#) algorithms for staffing, scheduling, rostering, routing.
- Integrate with [approximation](#) formulas and CTMC models.
- Better [models](#) for inbound traffic and other aspects.
- Currently, our work is driven by the interests of Bell Canada.

Papers: type [L'Ecuyer](#) in Google.