

# Stochastic Networks and Parameter Uncertainty

---

Assaf Zeevi\*

Graduate School of Business  
Columbia University

Stochastic Processing Networks Conference, August 2009

\* based on joint work with Mike Harrison Achal Bassamboo and Ramandeep  
Randhawa

## Motivation for this talk

---

Much of the work on stochastic processing networks:

- ▶ assumes model structure is *known a priori*, is *accurate* and *stationary*
  - parameters describing *system* and *environment* known and static
  - no model misspecification errors

## Motivation for this talk

---

### Much of the work on stochastic processing networks:

- ▶ assumes model structure is **known a priori**, is **accurate** and **stationary**
  - parameters describing **system** and **environment** known and static
  - no model misspecification errors

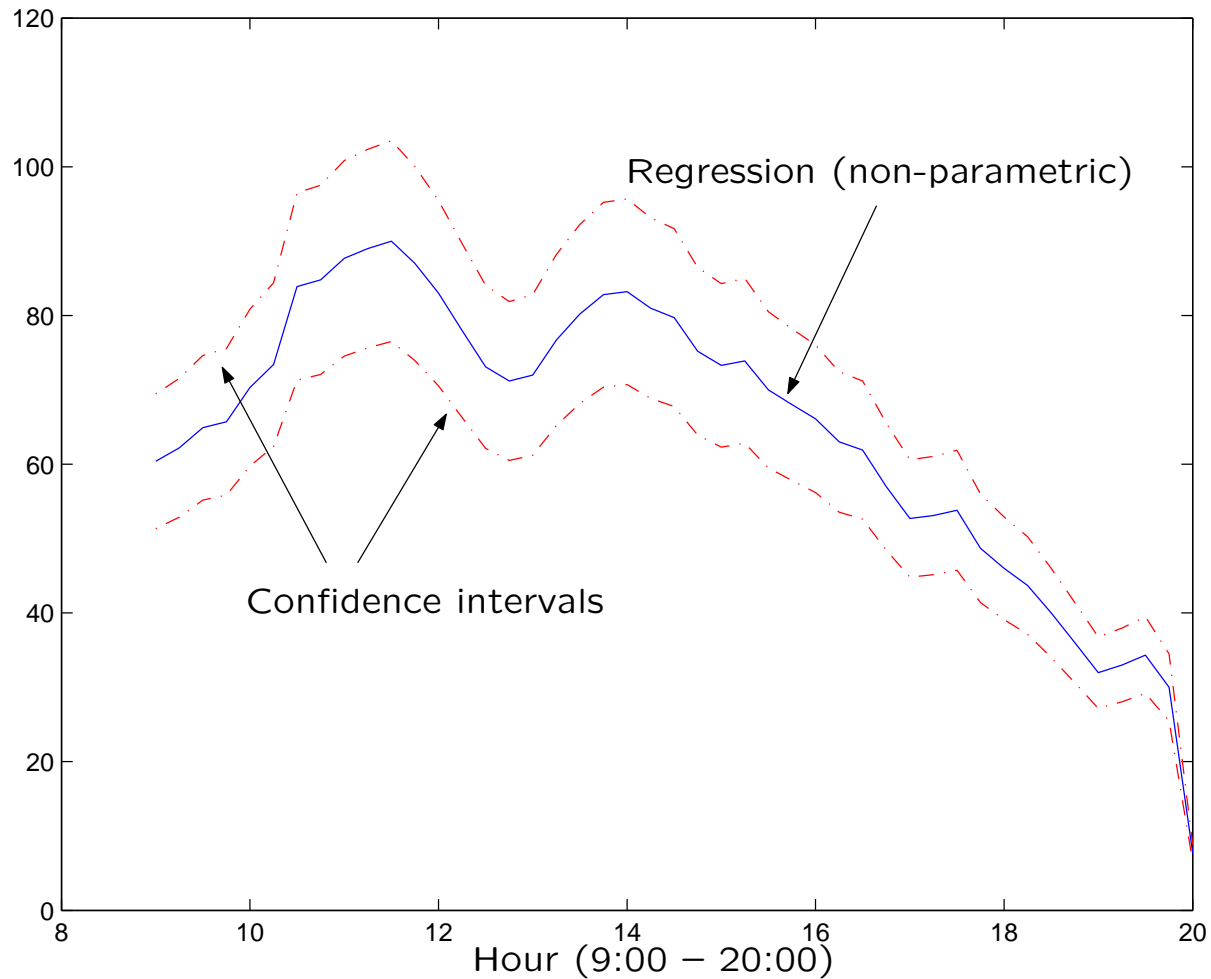
### In practice:

- ▶ model structure may be only **partially known**
- ▶ model primitives need to be **inferred**
  - from historical data
  - in on-line manner
- ▶ model may be **misspecified...**
  - both system model and environment
- ▶ environment may be **changing over time**

## Example 1: Design of global delivery centers

---

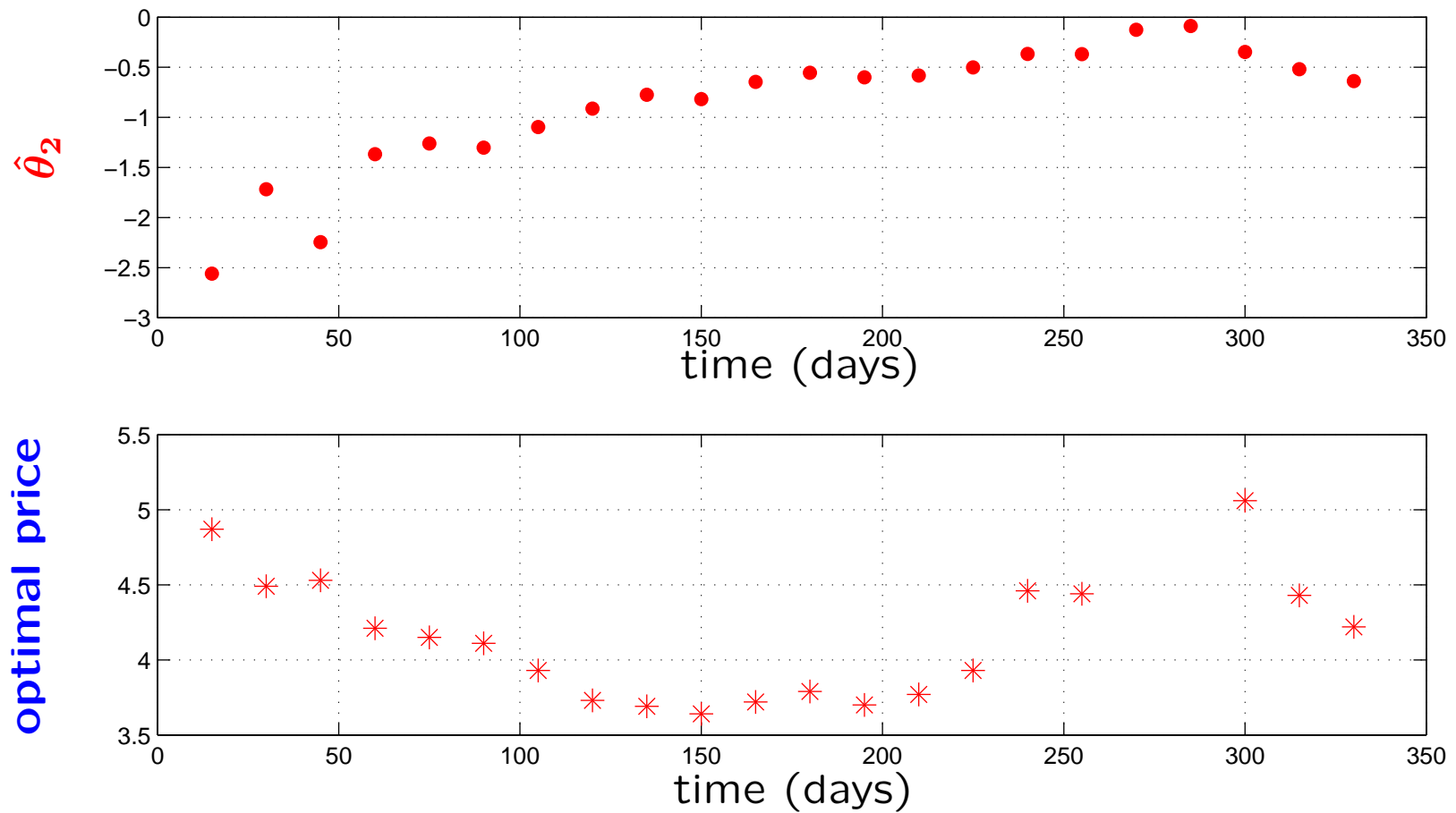
Arrival rate patterns in medium sized service center



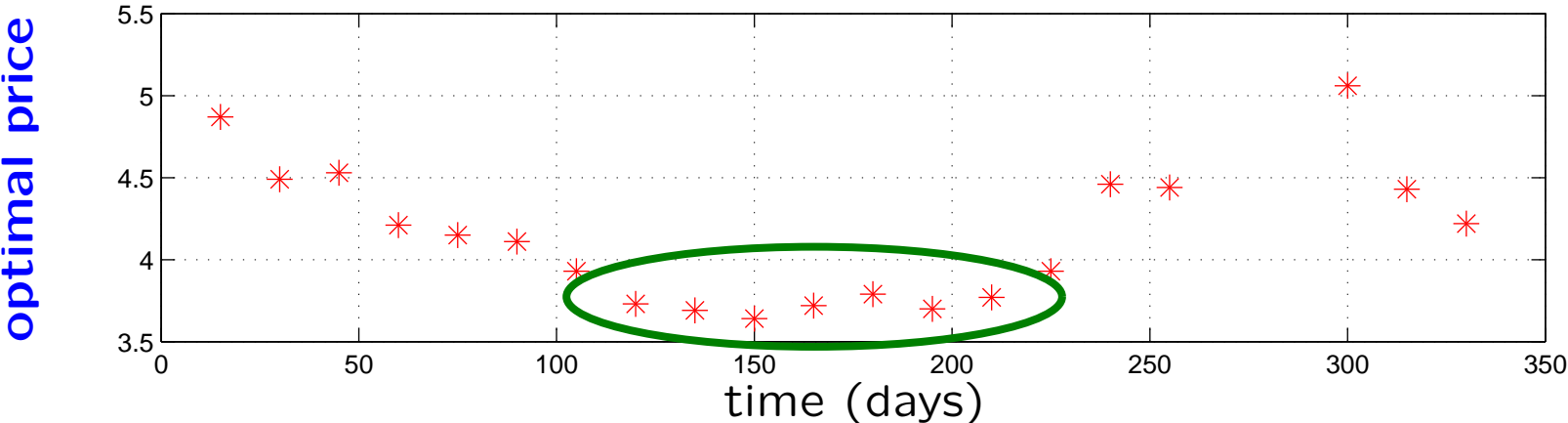
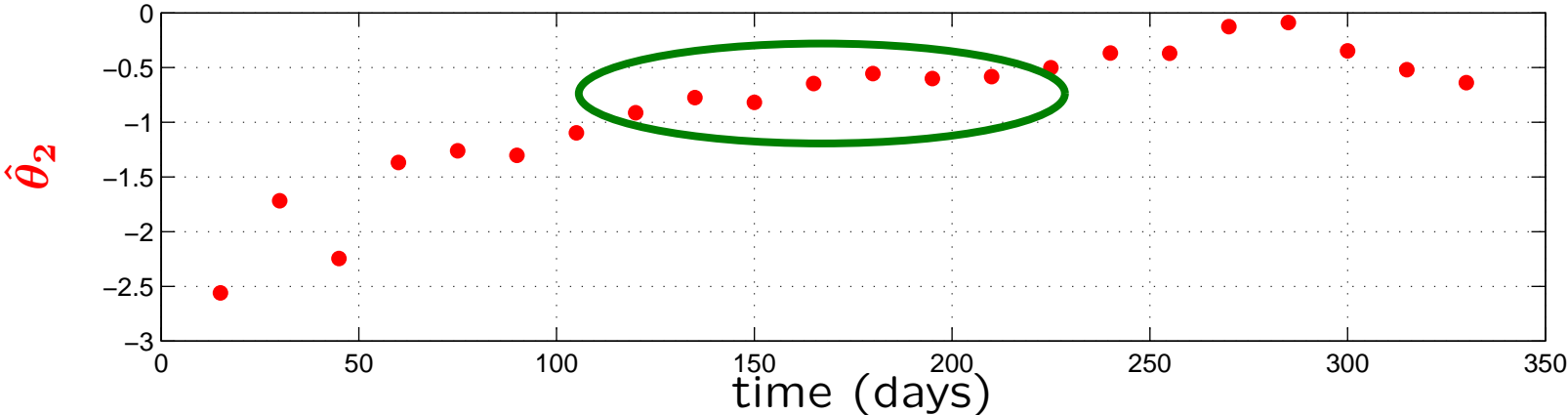
*how to deal with forecasting errors?*

## Example 2: Price engineering

---



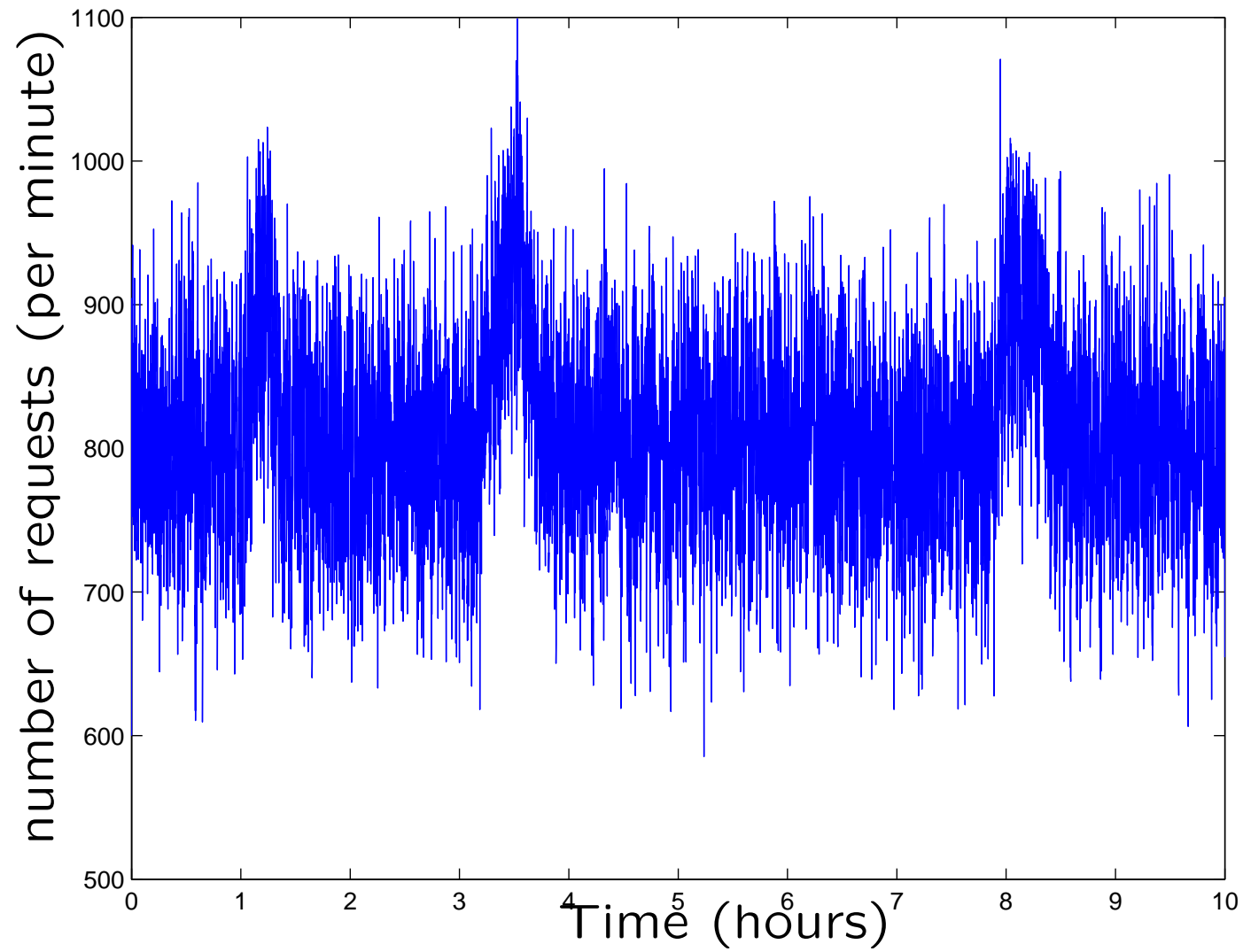
# Example 2: Price engineering



*how to deal with changing environment?*

## Example 3: Cloud computing

---



# What's in this talk

---

## Impact of parameter uncertainty on:

- ▶ static capacity / processing rate decisions
  - revisiting the square-root logic...
- ▶ model specification and calibration
  - estimation and testing
- ▶ dynamic control and resource allocation
  - revisiting the static planning problem...



## Parameter uncertainty and capacity planning

---

---

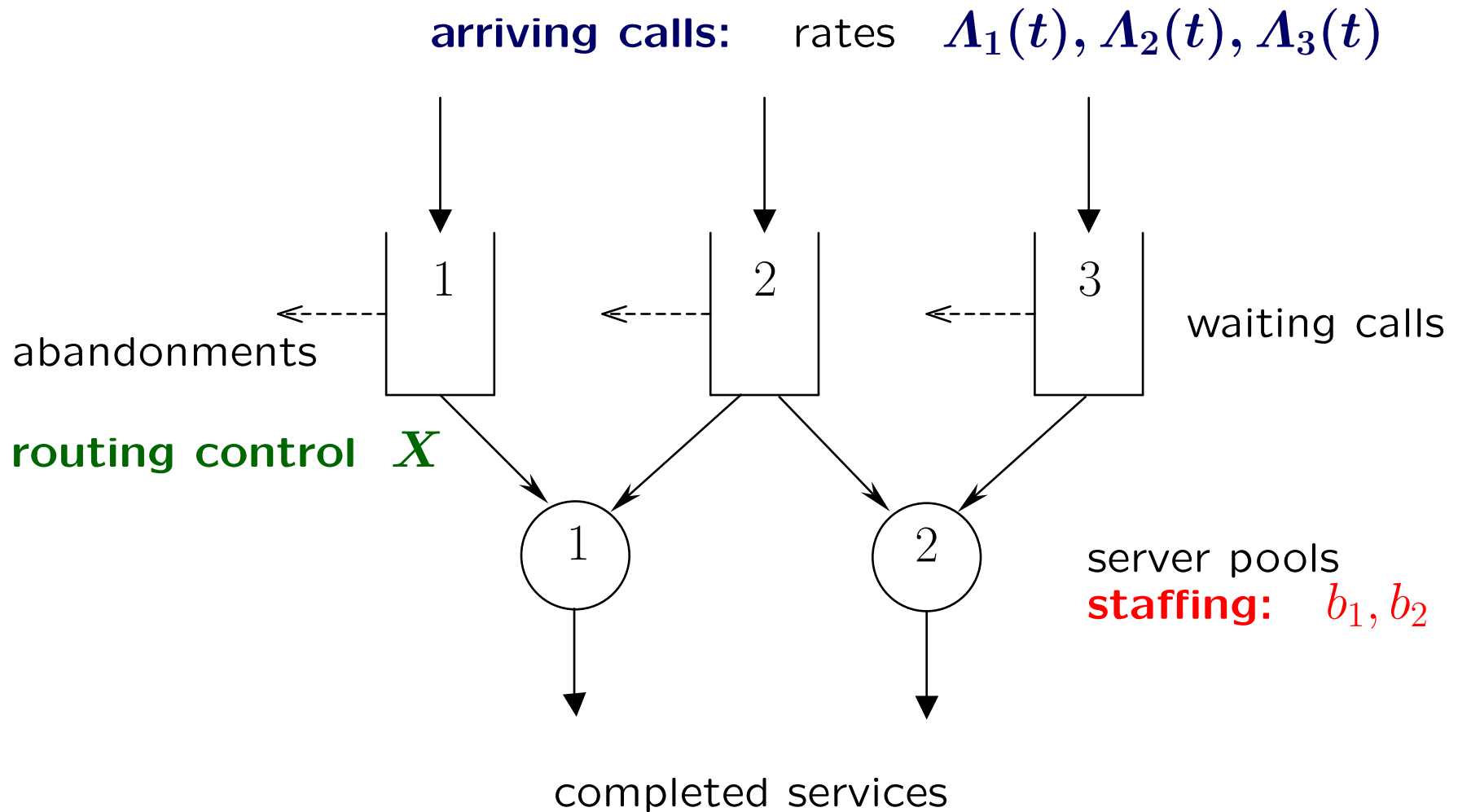
*Mean call arrivals 8 – 10AM in medium sized call center*

Day of Week	Mean no. of arriving calls	CV [empirical] (%)	CV [Poisson] (%)
Mon	943	26.5	3.3
Tue	824	22.3	3.5
Wed	807	26.5	3.5
Thu	778	28.5	3.6
Fri	767	33.5	3.6
Sat	293	61.8	5.8
Sun	139	148.1	8.5

- ▶ CV [empirical] = coefficient of variation (in %)
- ▶ CV [Poisson] = calculated *assuming* arrival process Poisson

# Parallel server network

---



- ▶ arrival process = doubly stochastic with rate  $\Lambda_1(t)$

## System dynamics

---

$$N_i(t) = \left[ \begin{array}{c} \text{Arrivals} \\ \end{array} \right] - \left[ \begin{array}{c} \text{Completed Services} \\ \end{array} \right] - \left[ \begin{array}{c} \text{Abandonments} \\ \end{array} \right]$$

- ▶  $N$  : **headcount process**

$N_i(t)$  = # of class  $i$  customers present at time  $t$

- ▶  $Q$  : **queue length process**

$Q_i(t)$  = # of class  $i$  customers not being served at time  $t$

- ▶  $\mathbf{X}$  : **dynamic control** [  $(R\mathbf{X})_i$  = rate of service in class  $i$  ]

$\mathbf{X}_j(t)$  = # of servers allocated to activity  $j$

- ▶  $b$  : **staffing vector**

- ▶  $(\mathbf{X}, N, Q)$  satisfy

$$A\mathbf{X}(t) \leq \mathbf{b}, \quad Q(t) = N(t) - B\mathbf{X}(t) \geq 0, \quad N(t) \geq 0, \quad \mathbf{X}(t) \geq 0$$

## System dynamics

---

$$N_i(t) = \left[ \begin{array}{c} \text{Arrivals} \\ \text{rate: } \Lambda_i(t) \end{array} \right] - \left[ \begin{array}{c} \text{Completed Services} \\ \text{rate: } (R\mathbf{X})_i(t) \end{array} \right] - \left[ \begin{array}{c} \text{Abandonments} \\ \text{rate: } \gamma_i Q_i(t) \end{array} \right]$$

- ▶  $N$  : **headcount process**

$N_i(t) = \#$  of class  $i$  customers present at time  $t$

- ▶  $Q$  : **queue length process**

$Q_i(t) = \#$  of class  $i$  customers not being served at time  $t$

- ▶  $\mathbf{X}$  : **dynamic control** [  $(R\mathbf{X})_i =$  rate of service in class  $i$  ]

$\mathbf{X}_j(t) = \#$  of servers allocated to activity  $j$

- ▶  $b$  : **staffing vector**

- ▶  $(\mathbf{X}, N, Q)$  satisfy

$$A\mathbf{X}(t) \leq \mathbf{b}, \quad Q(t) = N(t) - B\mathbf{X}(t) \geq 0, \quad N(t) \geq 0, \quad \mathbf{X}(t) \geq 0$$

# Design and control objectives

---

---

minimize:  $c \cdot \mathbf{b} + p \cdot \mathbb{E} \left[ \int_0^T \gamma Q(s) ds \right]$   
 $\mathbb{E} [\# \text{ of abandonments across classes}]$

capacity costs

s.t. **admissible routing control**  $\mathbf{X}$  over  $[0, T]$

---

$\mathbf{b}$  = r-dim'l vector of staffing levels in agent pools

$c$  = personnel cost vector

$p$  = penalty cost vector

$Q(t)$  = vector of queuelengths at time  $t$  in class  $i$  [depends on routing...]

$\gamma$  = abandonment rate vector

$T$  = planning horizon over which staffing is held fixed

## Design and control objectives

---

---

minimize:  $c \cdot \mathbf{b} + p \cdot \mathbb{E} \left[ \int_0^T \gamma Q(s) ds \right]$

capacity costs  $\mathbb{E} [\# \text{ of abandonments across classes}]$

s.t. **admissible routing control**  $\mathbf{X}$  over  $[0, T]$

---

$\mathbf{b}$  = r-dim'l vector of staffing levels in agent pools

$c$  = personnel cost vector

$p$  = penalty cost vector

$Q(t)$  = vector of queuelengths at time  $t$  in class  $i$  [depends on routing...]

$\gamma$  = abandonment rate vector

$T$  = planning horizon over which staffing is held fixed

**Decision “variables”:** capacity vector  $\mathbf{b}$  and control  $\mathbf{X}$

## A simple single-class / single-pool example

---

---

*"Solve the simplest problem you don't know the answer to."*

– Mike Harrison

## A simple single-class / single-pool example

---

---

*"Solve the simplest problem you don't know the answer to."*

– Mike Harrison

- ▶ arrival process doubly stochastic w/ rate  $\lambda(t)$
- ▶ exponential services w/ rate  $\mu$
- ▶ exponential reneging w/ rate  $\gamma$
- ▶  $b$  statistically identical servers



## A simple single-class / single-pool example

---

---

*"Solve the simplest problem you don't know the answer to."*

– Mike Harrison

- ▶ arrival process doubly stochastic w/ rate  $\Lambda(t)$
- ▶ exponential services w/ rate  $\mu$
- ▶ exponential reneging w/ rate  $\gamma$
- ▶  $b$  statistically identical servers

**objective:** minimize

$$\Pi(\mathbf{b}) := c \cdot \mathbf{b} + p \mathbb{E} \left[ \int_0^T \gamma Q(s) ds \right]$$

- ▶  $\mathbf{b}^*$  = optimal capacity choice

## Mike's key observation

---

**if**  $\Lambda \gg \mu, \gamma$  and of reasonable magnitude

- e.g., 100's of calls/hour, processing/renegeing order of minutes

**then** expect

- “accurate” fluid approximation
- “short” relaxation times

## Mike's key observation

---

---

**if**  $\Lambda \gg \mu, \gamma$  and of reasonable magnitude

- e.g., 100's of calls/hour, processing/renegeing order of minutes

**then** expect

- “accurate” fluid approximation
- “short” relaxation times

▶  $\gamma Q(t) \approx (\Lambda(t) - b\mu)^+$  *pointwise stationary fluid model* (PSFM)

## Mike's key observation

---

---

**if**  $\Lambda \gg \mu, \gamma$  and of reasonable magnitude

- e.g., 100's of calls/hour, processing/renegeing order of minutes

**then** expect

- “accurate” fluid approximation
- “short” relaxation times

▶  $\gamma Q(t) \approx (\Lambda(t) - b\mu)^+$  *pointwise stationary fluid model* (PSFM)

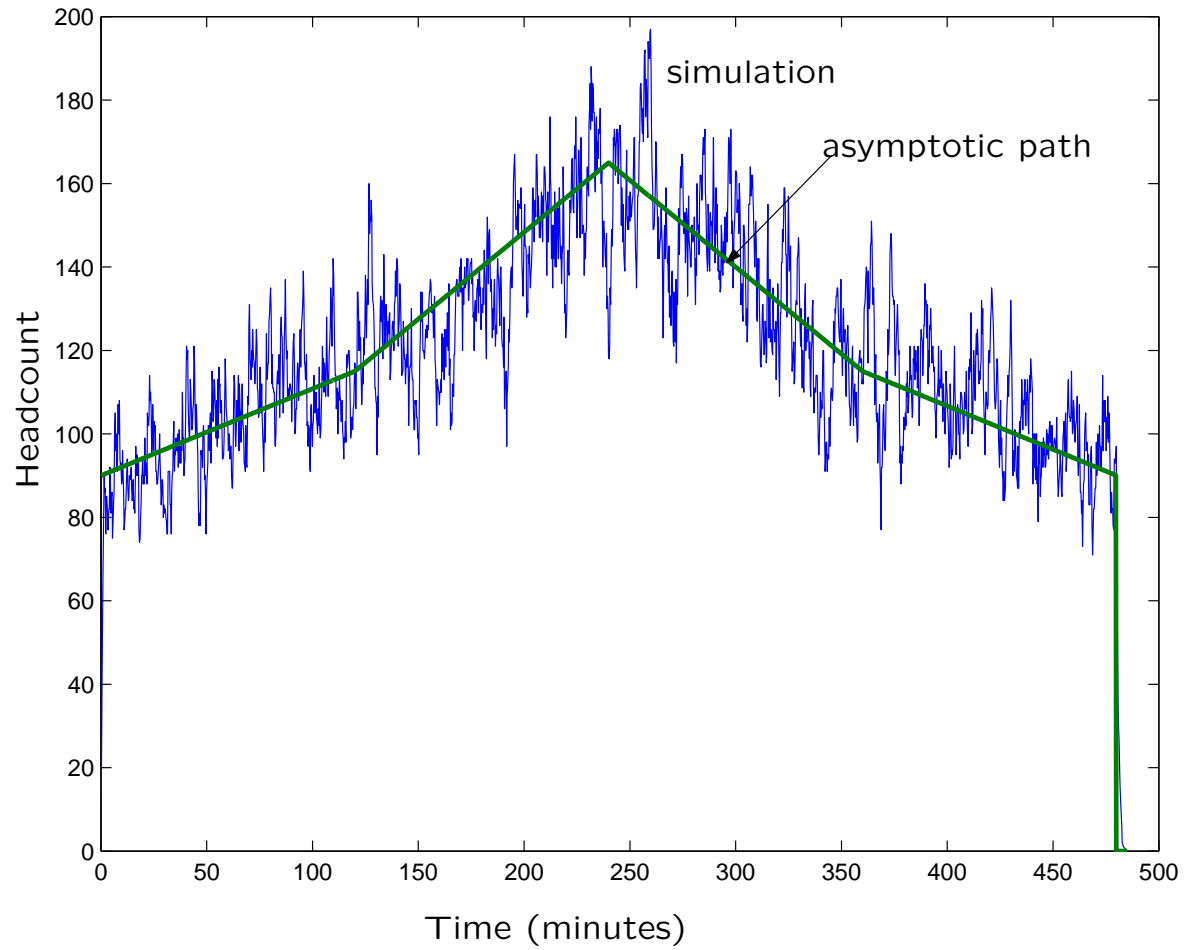
**More generally:** in each class  $i = 1, \dots, m$ , and all time  $t$

$$(*) \quad \gamma_i Q_i(t) \approx \Lambda_i(t) - (RX)_i(t)$$

abandonment rate      arrival rate      processing rate

# Picture proof

---



for real proof see [Bassamboo-Harrison-Z \(06a,06b\)](#)

## Consequence

---

original objective fn:

$$\Pi(\mathbf{b}) = c \cdot \mathbf{b} + p \cdot \mathbb{E}\left[\int_0^T \gamma Q(s) ds\right]$$

▶ optimal solution  $\mathbf{b}^*$

## Consequence

---

original objective fn:

$$\Pi(\mathbf{b}) = c \cdot \mathbf{b} + p \cdot \mathbb{E} \left[ \int_0^T \gamma Q(s) ds \right]$$

▶ optimal solution  $\mathbf{b}^*$

approximate objective fn:

$$\bar{\Pi}(\mathbf{b}) = c \cdot \mathbf{b} + p \cdot \mathbb{E} \left[ \int_0^T (\Lambda(s) - \mathbf{b}\mu)^+ ds \right]$$

▶ approximate solution  $\bar{\mathbf{b}}$

## Relation to traditional models

---

**if**  $\Lambda = \lambda$  (deterministic case)

**then** optimal solution takes form

$$b^* = \frac{\lambda}{\mu} + \beta \sqrt{\frac{\lambda}{\mu}}$$

▶ Erlang's square root rule...



## Relation to traditional models

---

**if**  $\Lambda = \lambda$  (deterministic case)

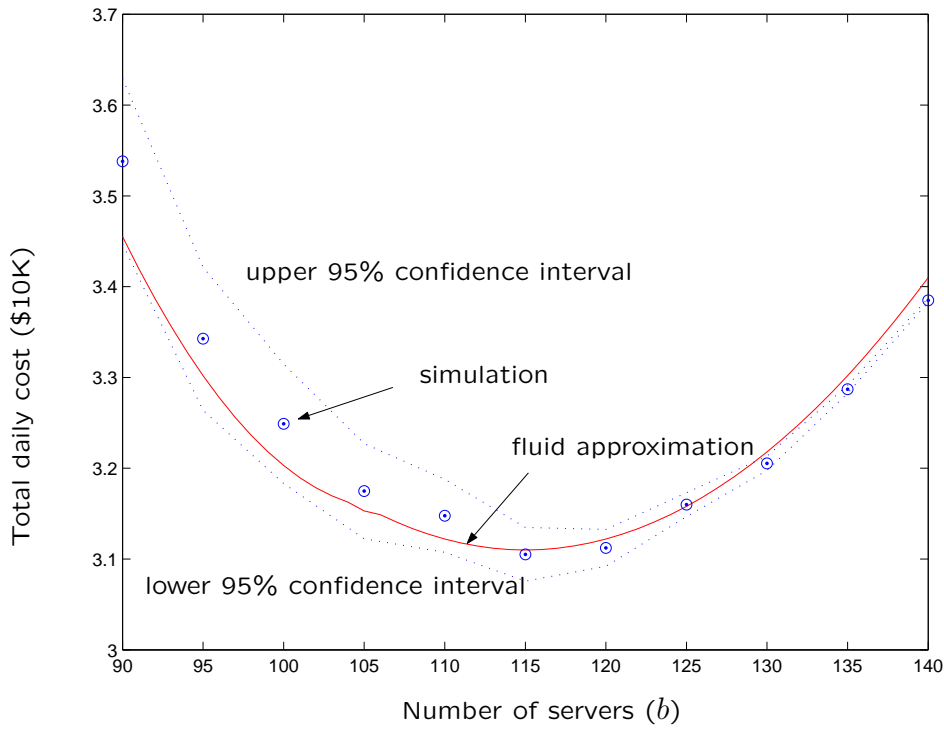
**then** optimal solution takes form

$$b^* = \frac{\lambda}{\mu} + \beta \sqrt{\frac{\lambda}{\mu}}$$

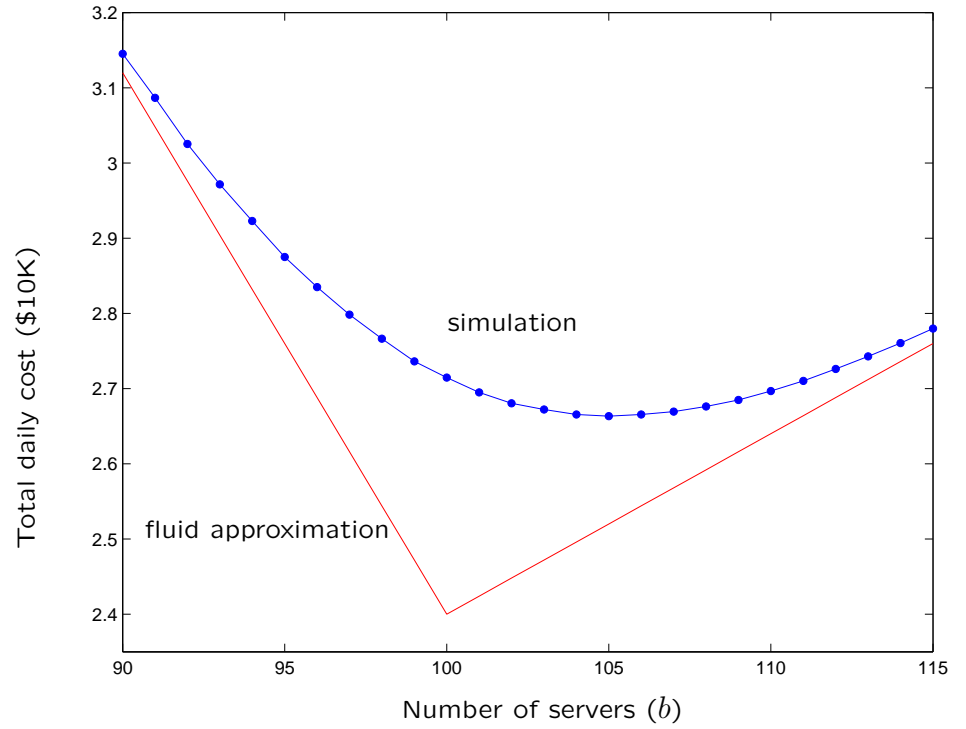
► Erlang's square root rule...

Arrival rate $\lambda$	Optimal solution		Prescription		Difference	
	$b^*$	$\Pi^*$	$\bar{b}$	$\Pi(\bar{b})$	$b^* - \bar{b}$	$\Pi(\bar{b}) - \Pi^*$
37.5	40	14.75	37	15.02	<b>3</b>	0.27
75	79	28.17	75	28.45	<b>4</b>	0.28
300	307	106.32	300	106.91	<b>7</b>	0.59

# In pictures...



Stochastic arrival rate



Deterministic arrival rate

## Relation to traditional models (cont'd)

---

---

- ▶ arrival rate  $\lambda$ 
  - time homogenous
  - *random* drawn from distribution  $F$

## Relation to traditional models (cont'd)

---

---

- ▶ arrival rate  $\lambda$ 
  - time homogenous
  - *random* drawn from distribution  $F$

**objective:** minimize

$$H(\mathbf{b}) := c \cdot \mathbf{b} + p \mathbf{E}[N - \mathbf{b}]^+$$

- ▶  $N$  = number of customers in system in steady-state

## Relation to traditional models (cont'd)

---

---

- ▶ arrival rate  $\Lambda$ 
  - time homogenous
  - *random* drawn from distribution  $F$

**objective:** minimize

$$\Pi(\mathbf{b}) := c \cdot \mathbf{b} + p \mathbf{E}[N - \mathbf{b}]^+$$

- ▶  $N$  = number of customers in system in steady-state

**PSFM approximation:** minimize

$$\bar{\Pi}(\mathbf{b}) := c \cdot \mathbf{b} + p \mathbf{E}[\Lambda - \mathbf{b}\mu]^+$$

- ▶ simple newsvendor problem with fractile solution

$$\bar{\mathbf{b}} = \frac{1}{\mu} \bar{F}^{-1} \left( \frac{c}{p\mu} \right)$$

## Accuracy of the newsvendor-based logic

*arrival rates constant and random (CV = 19.2%)*

Arrival rate distribution	Optimal solution		Prescription		Difference	
	$b^*$	$\Pi^*$	$\bar{b}$	$\Pi(\bar{b})$	$ b^* - \bar{b} $	$\Pi(\bar{b}) - \Pi^*$
U[25,50]	42	16.21	41	16.23	<b>1</b>	<b>0.02</b>
U[50,100]	83	31.47	83	31.47	<b>0</b>	<b>0</b>
U[200,400]	332	122.89	333	122.89	<b>1</b>	<b>0</b>

## Why is the prescription so accurate under uncertainty?

---

Consider simple case where  $\mu = \gamma$  [ just for purposes of intuition ]

- infinite server queue
- use normal approximation to Poisson...

$$\begin{aligned} \Pi(\mathbf{b}) &= c \cdot \mathbf{b} + p \mathbb{E} \left[ \int_0^T \gamma Q(s) ds \right] \\ &\approx \end{aligned}$$

## Why is the prescription so accurate under uncertainty?

---

Consider simple case where  $\mu = \gamma$  [ just for purposes of intuition ]

- infinite server queue
- use normal approximation to Poisson...

$$\begin{aligned} \Pi(\mathbf{b}) &= c \cdot \mathbf{b} + p \mathbb{E} \left[ \int_0^T \gamma Q(s) ds \right] \\ &\approx c \cdot \mathbf{b} + p\mu \mathbb{E}[\Lambda/\mu - \mathbf{b}]^+ + K \mathbb{E} \left[ \sqrt{\Lambda/\mu} \exp \left( -\frac{(\Lambda/\mu - \mathbf{b})^2}{2(\Lambda/\mu)} \right) \right] \end{aligned}$$



## Why is the prescription so accurate under uncertainty?

---

Consider simple case where  $\mu = \gamma$  [ just for purposes of intuition ]

- infinite server queue
- use normal approximation to Poisson...

$$\begin{aligned} \Pi(\mathbf{b}) &= c \cdot \mathbf{b} + p \mathbb{E} \left[ \int_0^T \gamma Q(s) ds \right] \\ &\approx c \cdot \mathbf{b} + p\mu \mathbb{E}[\Lambda/\mu - \mathbf{b}]^+ + K \mathbb{E} \left[ \sqrt{\Lambda/\mu} \exp \left( -\frac{(\Lambda/\mu - \mathbf{b})^2}{2(\Lambda/\mu)} \right) \right] \\ &= \text{approximate objective fn} + \text{approximation error} \end{aligned}$$

## Why is the prescription so accurate under uncertainty?

---

Consider simple case where  $\mu = \gamma$  [ just for purposes of intuition ]

- infinite server queue
- use normal approximation to Poisson...

$$\begin{aligned} \Pi(\mathbf{b}) &= c \cdot \mathbf{b} + p \mathbb{E} \left[ \int_0^T \gamma Q(s) ds \right] \\ &\approx c \cdot \mathbf{b} + p\mu \mathbb{E}[\Lambda/\mu - \mathbf{b}]^+ + K \mathbb{E} \left[ \sqrt{\Lambda/\mu} \exp \left( -\frac{(\Lambda/\mu - \mathbf{b})^2}{2(\Lambda/\mu)} \right) \right] \\ &= \text{approximate objective fn} + \text{approximation error} \end{aligned}$$

This suggests that *performance gap is bounded...*

- ▶ performance gap  $\Delta = \Pi(\bar{\mathbf{b}}) - \Pi^*$  is independent of scale of system...

## Rigorous foundations

---

Put  $\mathbb{E}\Lambda = n$  and  $CV^n =$  coefficient of variation

**Thm.** [ Bassamboo-Randhawa-Z (09) ]

▶ *Uncertainty-driven regime:* if  $CV^n \gg 1/\sqrt{n}$ , then

$$\Pi(\bar{\mathbf{b}}^n) = \Pi_* + \mathcal{O}(1/CV^n).$$

▶ *Variability-driven regime:* if  $CV^n \ll 1/\sqrt{n}$ , then

$$\Pi(\bar{\mathbf{b}}^n) = \Pi_* + \mathcal{O}(\sqrt{n}).$$

## Rigorous foundations

---

Put  $\mathbb{E}\Lambda = n$  and  $CV^n =$  coefficient of variation

**Thm.** [ Bassamboo-Randhawa-Z (09) ]

▶ *Uncertainty-driven regime:* if  $CV^n \gg 1/\sqrt{n}$ , then

$$\Pi(\bar{\mathbf{b}}^n) = \Pi_* + \mathcal{O}(1/CV^n).$$

▶ *Variability-driven regime:* if  $CV^n \ll 1/\sqrt{n}$ , then

$$\Pi(\bar{\mathbf{b}}^n) = \Pi_* + \mathcal{O}(\sqrt{n}).$$

**Cor 1.** *If CV bounded away from 0 then prescription is  $\mathcal{O}(1)$ -optimal*

## Rigorous foundations

---

Put  $\mathbb{E}\Lambda = n$  and  $CV^n =$  coefficient of variation

**Thm.** [ Bassamboo-Randhawa-Z (09) ]

▶ *Uncertainty-driven regime:* if  $CV^n \gg 1/\sqrt{n}$ , then

$$\Pi(\bar{\mathbf{b}}^n) = \Pi_* + \mathcal{O}(1/CV^n).$$

▶ *Variability-driven regime:* if  $CV^n \ll 1/\sqrt{n}$ , then

$$\Pi(\bar{\mathbf{b}}^n) = \Pi_* + \mathcal{O}(\sqrt{n}).$$

**Cor 1.** If  $CV$  bounded away from 0 then prescription is  $\mathcal{O}(1)$ -optimal

**Cor 2.** Performance of  $\bar{\mathbf{b}}^n$  is not sensitive to  $\mathcal{O}(\sqrt{n})$  perturbations

## Inference and model calibration

---

**problem:** previous slides assume distribution of arrival rate is known...

## Inference and model calibration

---

**problem:** previous slides assume distribution of arrival rate is known...

**possible approach:** [ Bassamboo- Z (2009) ]

- ▶ estimate arrival rate distribution  $F_n$  [  $n = \text{“sample size”}$  ]
- ▶ form empirical (approximate) objective fn  $\bar{II}_n(\cdot)$
- ▶ compute  $\bar{b}_n$  [ estimator of  $\bar{b}$  ]
- ▶ evaluate performance of estimator  $\mathbb{E}II(\bar{b}_n)$

## Inference and model calibration

---

**problem:** previous slides assume distribution of arrival rate is known...

**possible approach:** [ Bassamboo- Z (2009) ]

- ▶ estimate arrival rate distribution  $F_n$  [  $n = \text{“sample size”}$  ]
- ▶ form empirical (approximate) objective fn  $\bar{\Pi}_n(\cdot)$
- ▶ compute  $\bar{b}_n$  [ estimator of  $\bar{b}$  ]
- ▶ evaluate performance of estimator  $\mathbb{E}II(\bar{b}_n)$

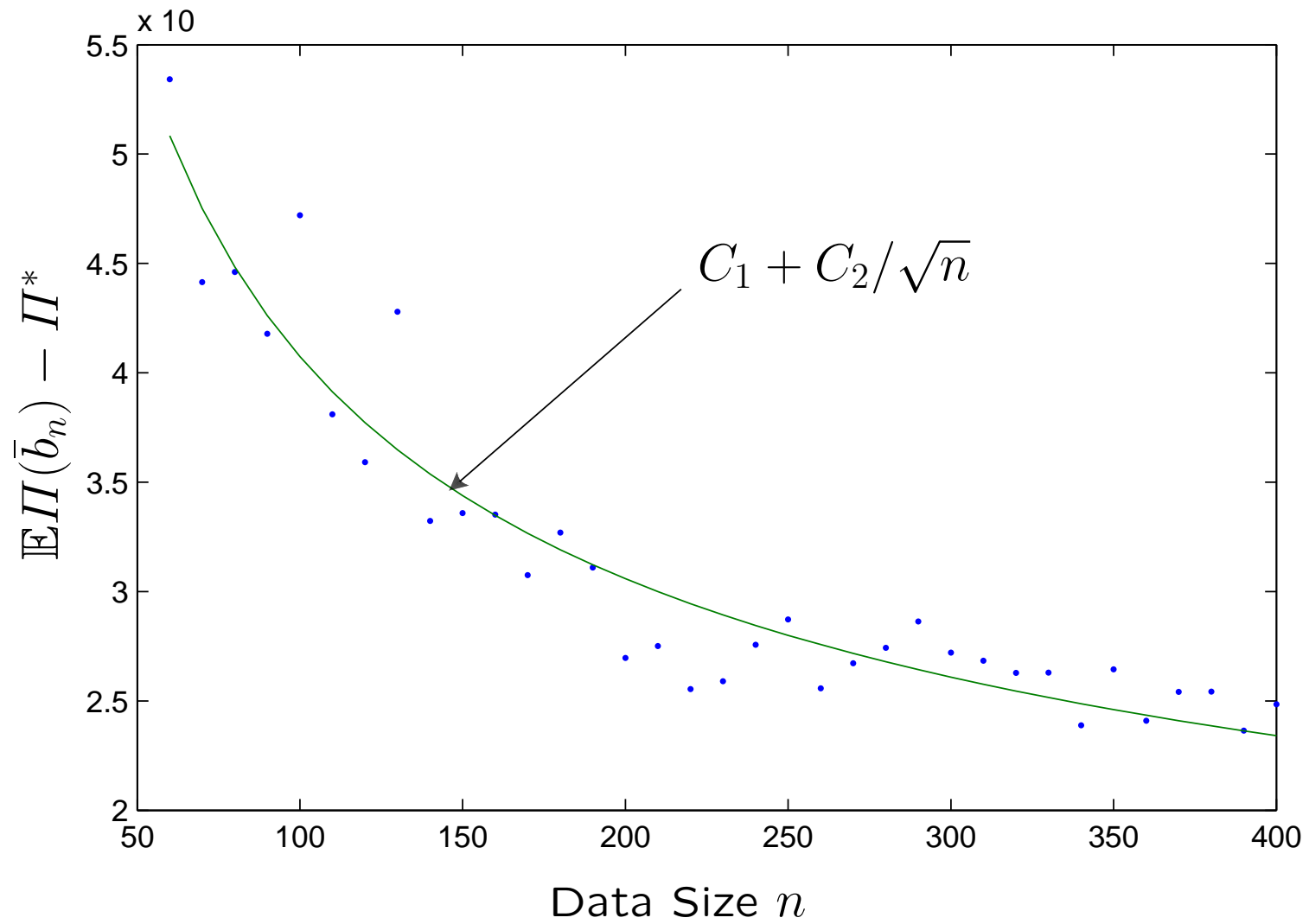
**key ideas in analysis:**

- ▶ need  $\bar{\Pi}_n(\cdot)$  to be amenable to  $M$ -estimation theory...
  - e.g., Lipschitz fn guarantees finite bracketing entropy
- ▶ use Talagrand's bounds to establish  $1/\sqrt{n}$  accuracy
  - $\mathbb{E}II(\bar{b}_n) - II^* = C/\sqrt{n} + \text{approximation error}$
  - interaction between approximation bound and estimation bound...



## Picture proof...

---



## Takeaway messages

---

### Parameter uncertainty:

- ▶ creates insensitivity in the objective fn
- ▶ makes it easier to achieve near optimal performance
  - simple capacity planning fluid problem
  - very simple control rules
- ▶ estimate/calibrate model
  - off line estimation [ capacity planning ]
  - real-time estimation [ control ]