

Data-Based Science for Service Engineering and Management

or: Empirical Adventures
in Call-Centers and Hospitals

Avi Mandelbaum

Technion, Haifa, Israel

<http://ie.technion.ac.il/serveng>

1

Research Partners

- **Students:**
Aldor*, Baron*, Carmeli, Feldman*, Garnett*, Gurvich*, Huang, Khudiakov*, Maman*, Marmor*, Reich, Rosenshmidt*, Shaikhet*, Senderovic, Tseytlin*, Yom-Tov*, Yuviler, Zaled, Zeltyn*, Zychlinski, Zohar*, Zviran*, ...
- **Theory:**
Armory, Atar, Gurvich, Jelenkovic, Kaspi, Massey, Momcilovic, Reiman, Shimkin, Stolyar, Wasserkrug, Whitt, Zeltyn, ...
- **Industry:**
Mizrahi Bank (A. Cohen, U. Yonissi), Rambam Hospital (R. Beyar, S. Israelit, S. Tzafrir), IBM Research (OCR Project), Hapoalim Bank (G. Maklef, T. Shlasky), Pelephone Cellular, ...
- **Technion SEE Center / Laboratory:**
Feigin; Trofimov, Nadjharov, Gavako, Kutsy; Liberman, Koren, Plonsky, Senderovic; Research Assistants, ...
- **Empirical/Statistical Analysis:**
Brown, Gans, Zhao; Shen; Ritov, Goldberg; Gurvich, Huang, Liberman; Armory, Marmor, Tseytlin, Yom-Tov; Zeltyn, Nardi, Gorfine, ...

2

History, Resources (Downloadable)

- Math. + C.S. + Stat. + O.R. + Mgt. ⇒ **IE** (≥ 1990)
- **Teaching: "Service-Engineering" Course** (≥ 1995):
<http://ie.technion.ac.il/serveng> - website
http://ie.technion.ac.il/serveng/References/teaching_paper.pdf
- **Call-Centers Research** (≥ 2000)
e.g. <Call Centers> in Google-Scholar
- **Healthcare Research** (≥ 2005)
e.g. **OCR Project**: IBM + Rambam Hospital + Technion
- **The Technion SEE Center** (≥ 2007)

3

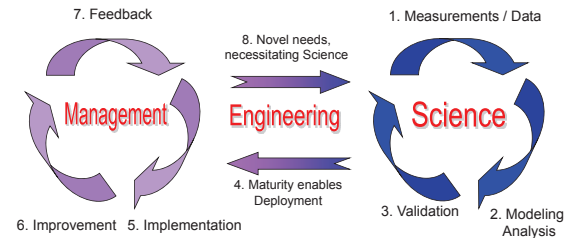
The Case for Service Science / Engineering

- **Service Science / Engineering** (vs. Management) are emerging **Academic Disciplines**. For example, universities (world-wide), IBM (SSME, a la Computer-Science), USA NSF (SEE), Germany IAO (ServEng), ...
- Models that explain **fundamental phenomena**, which are **common** across applications:
 - **Call Centers**
 - **Hospitals**
 - **Transportation**
 - Justice, Fast Food, Police, Internet, ...
- **Simple models** at the Service of **Complex Realities** (Human)
Note: Simple yet rooted in **deep analysis**.
- Mostly **What Can Be Done** vs. **How To**

4

Title: Expands the Scientific Paradigm

Physics, Biology, ... : **Measure, Model, Experiment, Validate, Refine.**
Human-complexity triggered above in Transportation, Economics.
Starting with **Data**, expand to:

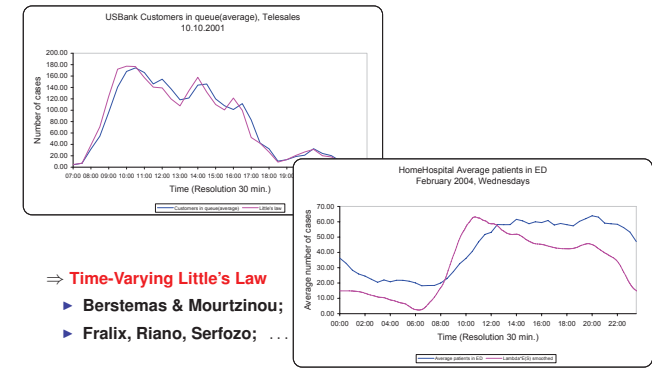


e.g. Validate, refute or discover **congestion laws** (Little, PASTA, SSC, ?, ?, ...), in call centers and hospitals

5

Little's Law: Call Center & Emergency Department

Time-Gap: # in System lags behind **Piecewise-Little** ($L = \lambda \times W$)



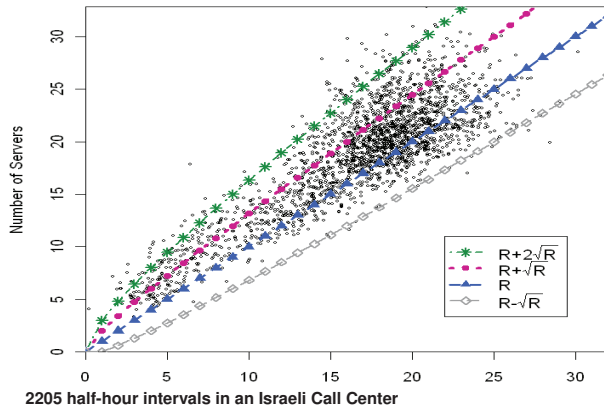
⇒ **Time-Varying Little's Law**

- Berstemas & Mourtzinou;
- Fralix, Riano, Serfozo; ...

6

QED Call Center: Staffing (N) vs. Offered-Load (R)

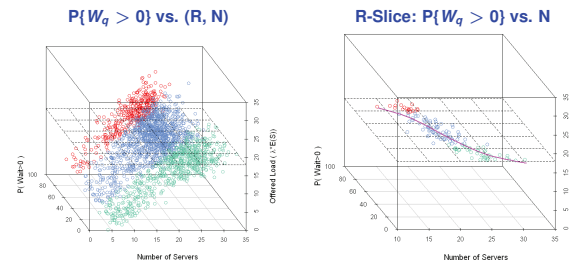
IL Telecom; June-September, 2004; w/ Nardi, Plonski, Zeltyn



7

QED Call Center: Performance

Large Israeli Bank



3 Operational Regimes:

- **QD:** ≤ 25%
- **QED:** 25% – 75%
- **ED:** ≥ 75%

8

Operational Regimes: Scaling, Performance, w/ I. Gurvich & J. Huang

Erlang-A n fixed	Conventional scaling			MS scaling			NDS scaling			
	Sub	Critical	Super	QD	QED	ED	ED+QED	Sub	Critical	Super
Offered load per server	$\frac{\lambda}{n} < 1$	$1 - \frac{\lambda}{n} \approx 1$	$\frac{\lambda}{n} > 1$	$\frac{\lambda}{n}$	$1 - \frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n} - \beta \sqrt{\frac{\lambda}{n}}$	$\frac{\lambda}{n}$	$1 - \frac{\lambda}{n}$	$\frac{\lambda}{n}$
Arrival rate λ	$\frac{\lambda}{n}$	$\mu - \frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\mu - \beta \mu \sqrt{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n} - \beta \mu \sqrt{\frac{\lambda}{n}}$	$\frac{\lambda}{n}$	$\mu - \beta \mu$	$\frac{\lambda}{n}$
Number of servers	1	n			n			1	n	
Time-scale	n	n			n			n	n	
Abandonment rate	θ/n	n			θ			θ/n	n	
Staffing level	$\frac{\lambda}{n}(1 + \delta)$	$\frac{\lambda}{n}(1 + \frac{\delta}{\sqrt{n}})$	$\frac{\lambda}{n}(1 + \gamma)$	$\frac{\lambda}{n}(1 + \delta)$	$\frac{\lambda}{n}(1 + \frac{\delta}{\sqrt{n}})$	$\frac{\lambda}{n}(1 + \gamma)$	$\frac{\lambda}{n}(1 + \gamma) + \beta \sqrt{\frac{\lambda}{n}}$	$\frac{\lambda}{n}(1 + \delta)$	$\frac{\lambda}{n}(1 + \gamma)$	$\frac{\lambda}{n}(1 + \gamma)$
Utilization	$\frac{\lambda}{n}$	$1 - \frac{\lambda}{n}$	1	$\frac{\lambda}{n}$	$1 - \frac{\lambda}{n}$	1	1	$\frac{\lambda}{n}$	$1 - \frac{\lambda}{n}$	1
$E(Q)$	$\frac{\lambda}{n}$	$\sqrt{\frac{\lambda}{n}}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\sqrt{\frac{\lambda}{n}}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\sqrt{\frac{\lambda}{n}}$	$\frac{\lambda}{n}$
$P(A_k)$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$	$\frac{1}{k!} \frac{\lambda^k}{k!} e^{-\lambda/n}$
$P(W_q > 0)$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$
$\alpha_1 \in (0, 1)$	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1
$P(W_q > T)$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$	$\alpha_1 e^{-\frac{\lambda}{n}}$
Congestion	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$

- $\delta > 0, \gamma \in (0, 1)$ and $\beta \in (-\infty, \infty)$;
- $QD = \frac{\lambda}{n} < 1$;
- $ED = \frac{\lambda}{n} > 1$; $G(T) = \gamma$;
- $QED = \frac{\lambda}{n} \approx 1$; $\frac{\lambda}{n} = 1 - \beta \sqrt{\frac{\lambda}{n}}$ and $\beta = \frac{\lambda}{n}$;
- $ED+QED = \frac{\lambda}{n} \approx 1$; $G(T) = \gamma$;
- Conventional: critical: $P(W > T) = \frac{\lambda}{n}$; super: $P(W > T) = \frac{\lambda}{n}$; NDS: super: $P(W > T) = \frac{\lambda}{n}$;

9

Prerequisite I: Data

Averages Prevalent (and could be useful / interesting).

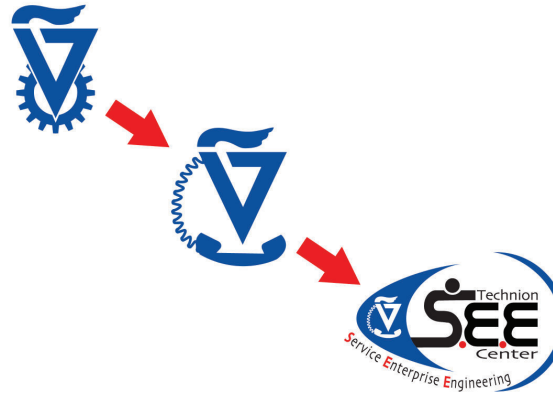
But I need data at the level of the **Individual Transaction**:
For each service transaction (during a phone-service in a call center, or a patient's visit in a hospital, or browsing in a website, or ...), its **operational history** = **time-stamps of events**.

Sources: **"Service-floor"** (vs. Industry-level, Surveys, ...)

- Administrative (Court, via "paper analysis")
- Face-to-Face (Bank, via bar-code readers)
- Telephone (Call Centers, via ACD / CTI, IVR/VRU)
- Hospitals (Emergency Departments, ...)
- Expanding:
 - Hospitals, via **RFID**
 - Operational + Financial + Contents (Marketing, Clinical)
 - Internet, Chat (multi-media)

10

Pause for a Commercial: The Technion SEE Center



11

Technion SEE = Service Enterprise Engineering

SEELab: Data-repositories for research and teaching

- For example:
 - Bank Anonymous: **1 years, 350K calls by 15 agents** - in 2000. **Brown, Gans, Sakov, Shen, Zeltyn, Zhao** (JASA), paved the way for:
 - U.S. Bank: **2.5 years, 220M calls, 40M by 1000 agents**.
 - Israeli Cellular: **2.5 years, 110M calls, 25M calls by 750 agents**.
 - Israeli Bank: **from January 2010, daily-deposit** at a SEESafe.
 - Israeli Hospital: **4 years, 1000 beds; 8 ED's- Sinreich's data**.

SEESat: Environment for graphical EDA in real-time

- Universal Design, Internet Access, Real-Time Response.

SEEServer: Free for academic use

Register, then access (presently) U.S. Bank and Bank Anonymous.

Visitor: run mstsc, seeserver.iem.technion.ac.il; Self-Tutorial

12

Tutorial Cover; State-Space Collapse from Tutorial

4 overheads:

- Cover (make sure relevant to the lecture (e.g. APS, HKUST))
- Page 2 (again, make sure relevant to the lecture)
- Contents (with Stat-Space Collapse yellowed)
- The page with State-Space Collapse.

13

eg. RFID-Based Data: Mass Casualty Event (MCE)

Drill: Chemical MCE, Rambam Hospital, May 2010



Focus on **severely wounded** casualties (≈ 40 in drill)

Note: 20 observers support real-time control (helps validation)

14

Data Cleaning: MCE with RFID Support

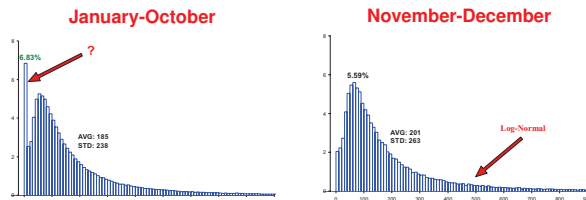
Data-base				Company report		comment
Asset id	order	Entry date	Exit date	Entry date	Exit date	
4	1	1:14:07 PM		1:14:00 PM		
6	1	12:02:02 PM	12:33:10 PM	12:02:00 PM	12:33:00 PM	
11	1	11:37:15 AM	12:40:17 PM	11:37:00 AM		exit is missing
10	1	12:23:32 PM	12:38:23 PM	12:23:00 PM		
12	1	12:12:47 PM	12:35:33 PM		12:35:00 PM	entry is missing
15	1	1:07:15 PM		1:07:00 PM		
16	1	11:18:19 AM		11:18:00 AM	11:31:00 AM	
17	1	1:03:31 PM	11:31:04 AM	1:03:00 PM		
18	1	1:07:54 PM		1:07:00 PM		
19	1	12:01:58 PM		12:01:00 PM		
20	1	11:37:21 AM	12:57:02 PM	11:37:00 AM	12:57:00 PM	
21	1	12:01:16 PM	12:37:16 PM	12:01:00 PM		
22	1	12:04:31 PM	12:20:40 PM			first customer is missing
22	2	12:27:37 PM		12:27:00 PM		
25	1	12:27:35 PM	1:07:28 PM	12:27:00 PM	1:07:00 PM	
27	1	12:06:53 PM		12:06:00 PM		exit time instead of entry time
28	1	11:21:34 AM	11:41:06 AM	11:21:00 AM	11:53:00 AM	
29	1	12:21:06 PM	12:54:29 PM	12:21:00 PM	12:54:00 PM	
31	1	11:40:54 AM	12:30:16 PM	11:40:00 AM	12:30:00 PM	
31	2	12:37:57 PM	12:54:51 PM	12:37:00 PM	12:54:00 PM	
32	1	11:27:11 AM	12:15:17 PM	11:27:00 AM	12:15:00 PM	
33	1	12:05:50 PM	12:13:12 PM	12:05:00 PM	12:15:00 PM	wrong exit time
35	1	11:31:48 AM	11:40:50 AM	11:31:00 AM	11:40:00 AM	
36	1	12:06:23 PM	12:29:30 PM	12:06:00 PM	12:29:00 PM	
37	1	11:31:50 AM	11:48:18 AM	11:31:00 AM	11:48:00 AM	
37	2	12:59:21 PM		12:59:00 PM		

Imagine **"Cleaning" 60,000+ customers per day** (call centers) !

15

Beyond Averages: The Human Factor

Histogram of Service-Time in a (Small Israeli) Bank, 1999

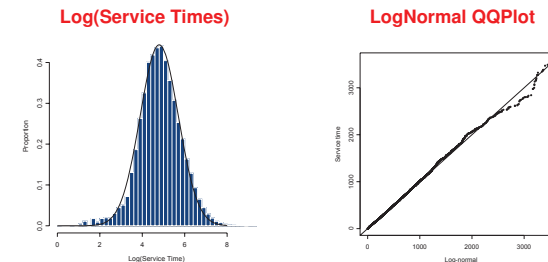


- 6.8% Short-Services:** Agents' "Abandon" (improve bonus, rest), (mis)lead by **incentives**
- Distributions** must be measured (in **seconds = natural scale**)
- LogNormal** service times common in call centers

16

Validating LogNormality of Service-Duration

Israeli Call Center, Nov-Dec, 1999

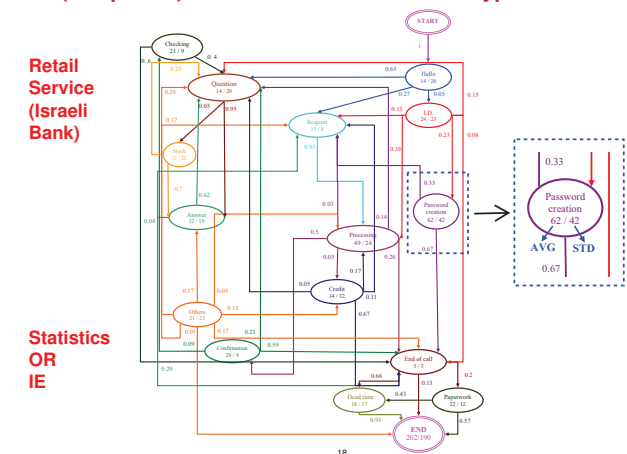


- Practically Important:** (mean, std)(log) characterization
- Theoretically Intriguing:** Why LogNormal? Naturally multiplicative but, in fact, also **Infinitely-Divisible** (Generalized Gamma-Convolutions)
- Simple-model of a complex-reality? The **Service Process**:

17

(Telephone) Service-Process = "Phase-Type" Model

Retail Service (Israeli Bank)



Statistics OR IE

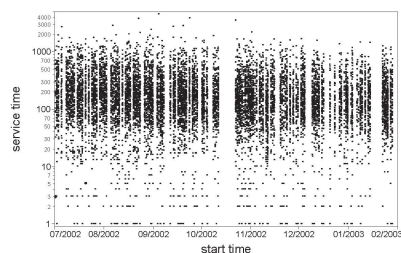
18

Individual Agents: Service-Duration, Variability

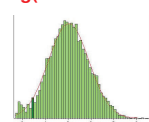
w/ Gans, Liu, Shen & Ye

Agent 14115

Service-Time Evolution: 6 month



Log(Service-Time)

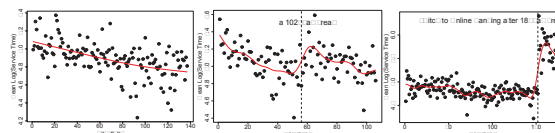


- ▶ **Learning:** Noticeable decreasing-trend in service-duration
- ▶ **LogNormal** Service-Duration, individually and collectively

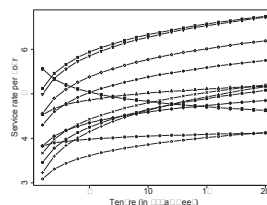
19

Individual Agents: Learning, Forgetting, Switching

Daily-Average Log(Service-Time), over 6 months
Agents 14115, 14128, 14136



Weakly Learning-Curves for 12 Homogeneous(?) Agents



Why Bother?

In large call centers:

+One Second to Service-Time implies +Millions in costs, annually

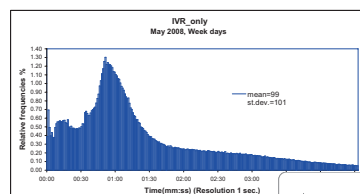
⇒ Time and "Motion" Studies (Classical IE with New-age IT)

- ▶ **Service-Process Model:** Customer-Agent Interaction
 - ▶ **Work Design** (w/ Khudiakov)
 - eg. **Cross-Selling:** higher profit vs. longer (costlier) services; Analysis yields (congestion-dependent) cross-selling protocols
 - ▶ **"Worker" Design** (w/ Gans, Liu, Shen & Ye)
 - eg. **Learning, Forgetting, ...** : Staffing & individual-performance prediction, in a heterogenous environment
- ▶ **IVR-Process Model:** Customer-Machine Interaction
 - 75% bank-services**, poor design, yet scarce research; Same approach, automatic (easier) data (w/ Yuviler)

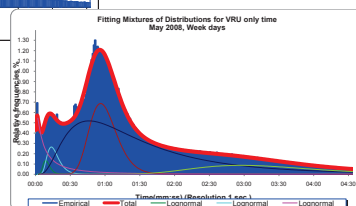
21

IVR-Time: Histograms

Israeli Bank: IVR/VRU Only, May 2008

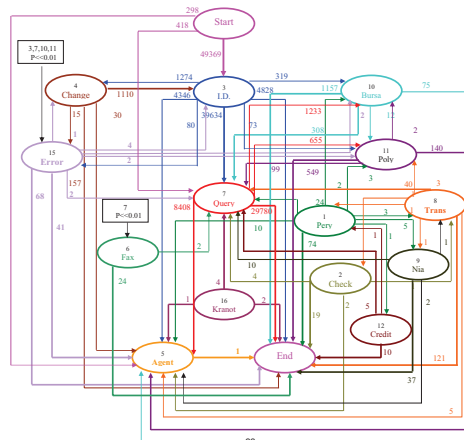


Mixture: 7 LogNormals



22

IVR-Process: "Phase-Type" Model



23

Started with Call Centers, Expanded to Hospitals

Call Centers - U.S. (Netherlands) Stat.

- ▶ \$200 – \$300 billion annual expenditures (0.5)
- ▶ 100,000 – 200,000 call centers (1500-2000)
- ▶ "Window" into the company, for better or worse
- ▶ Over 3 million agents = **2% – 4% workforce** (100K)

Healthcare - similar and unique challenges:

- ▶ Cost-figures far more staggering
- ▶ Risks much higher
- ▶ ED (initial focus) = hospital-window
- ▶ Over 3 million nurses

24

Call-Center Environment: Service Network



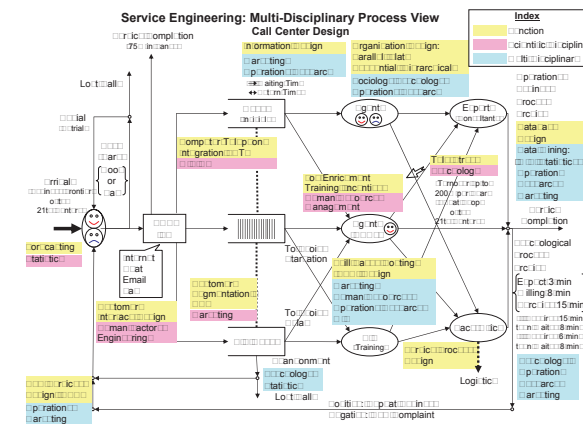
25

Call-Centers: "Sweat-Shops of the 21st Century"



26

Call-Center Network: Gallery of Models



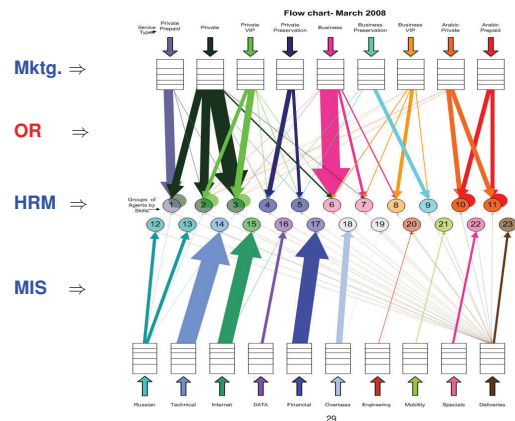
27

Call-Center Network: Gallery of Models

Add marks of topics to focus on

28

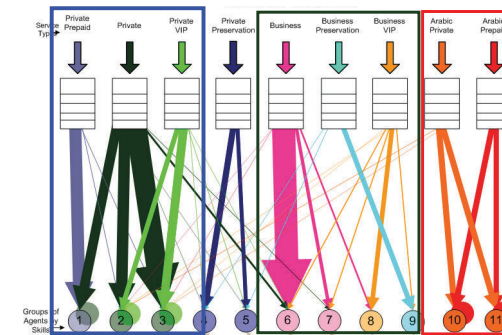
Skills-Based Routing in Call Centers EDA and OR, with I. Gurvich and P. Liberman



29

SBR Topologies: I; V, Reversed-V; N, X; W, M

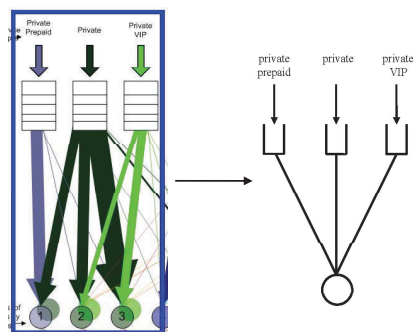
Israeli Cellular, March 2008



30

SBR: Class-Dependent Services

“Reduction” to V-Topology (Equivalent Brownian Control)

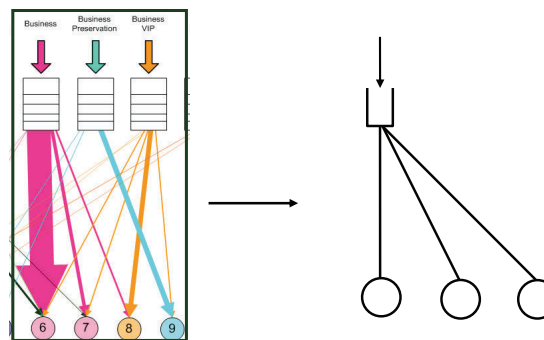


PhD's: Tezcan, Dai; Shaikhet, w/ Atar; Gurvich, Whitt

31

SBR: Pool-Dependent Services

“Reduction” to Reversed-V and I (Equivalent Brownian Control)

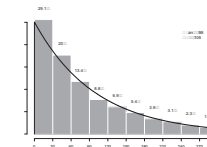


PhD's: Tezcan, Dai; Shaikhet, w/ Atar; Gurvich, Whitt

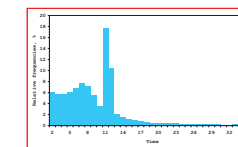
32

Waiting Times in a Call Center (Theory?)

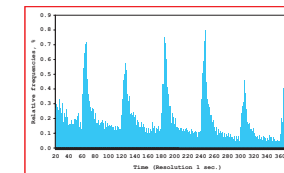
Exponential in Heavy-Traffic (min.)
Small Israeli Bank



Routing via Thresholds (sec.)
Large U.S. Bank



Scheduling Priorities (sec) (later: Hospital LOS, hr.)
Medium Israeli Bank



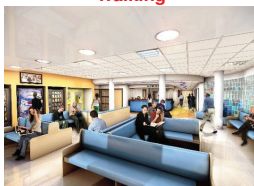
33

ER / ED Environment: Service Network

Acute (Internal, Trauma)



Walking



Multi-Trauma



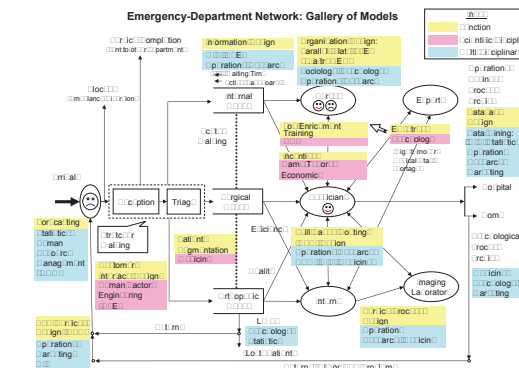
34

Queueing in a “Good” Beijing Hospital, at 6am



35

Emergency-Department Network: Gallery of Models



► Forecasting, Abandonment = LWBS, SBR ≈ Flow Control

36

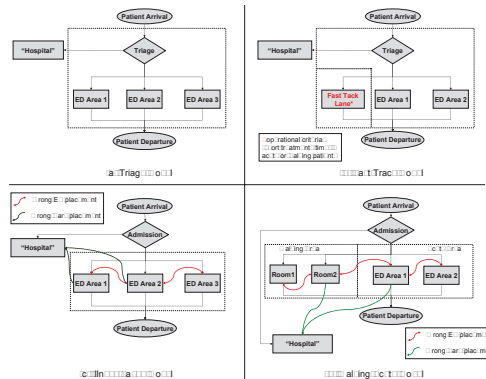
Emergency-Department Network: Gallery of Models

Add ED-to-IW routing

37

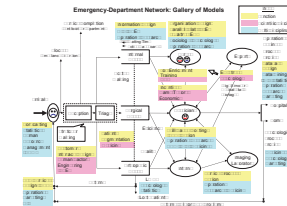
ED Design, with B. Golany, Y. Marmor, S. Israelit

Routing: Triage (Clinical), Fast-Track (Operational), ... (via DEA)
eg. Fast Track most suitable when elderly dominate



38

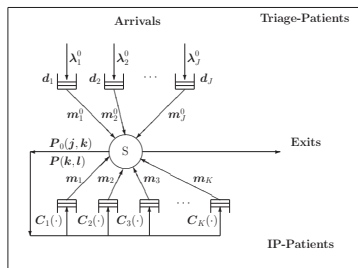
Emergency-Department Network: Flow Control



- ▶ Queueing-Science, w/ Armony, Marmor, Tseytlin, Yom-Tov
- ▶ Fair ED-to-IW Routing (Patients vs. Staff), w/ Momcilovic, Tseytlin
- ▶ Triage vs. In-Process / Release in EDs, w/ Carmeli, Huang, Shimkin
- ▶ Workload and Offered-Load in Fork-Join Networks, w/ Kaspi, Zaid
- ▶ Synchronization Control of Fork-Join Networks, w/ Atar, Zviran
- ▶ Staffing Time-Varying Q's with Re-Entrant Customers, w/ Yom-Tov

39

ED Patient Flow: The Physicians View



- ▶ Goal: Adhere to Triage-Constraints, then process/release In-Process Patients
- ▶ Model = Multi-class Q with Feedback: Min. convex congestion costs of IP-Patients, s.t. deadline constraints on Triage-Patients.
- ▶ Solution: In conventional heavy-traffic, asymptotic least-cost s.t. asymptotic compliance, via threshold (w/ B. Carmeli, J. Huang, S. Israelit, N. Shimkin; as in Plambeck, Harrison, Kumar, who applied admission control).

40

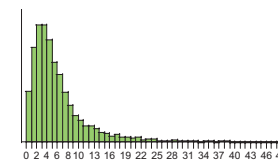
Operational Fairness

1. "Punishing" fast wards in ED-to-IW Routing:
 - ▶ Parallel IWs: similar clinically, differ operationally
 - ▶ Problem: Short Length-of-Stay goes hand in hand with high bed-occupancy, bed-turnover, yet clinically apt: **unfair!**
 - ▶ Solution: Both nurses and managers content, w/ P. Momcilovic and Y. Tseytlin (3 time-scales: hour, day, week; "compare" with call-centers SBR)
 2. Balancing Load across Maternity Wards:
 - ▶ 2 Maternity Wards: 1 = pre-birth, 2 = post-birth complications
 - ▶ Problem: Nurses think the "others-work-less": **unfair!**
 - ▶ Goal: Balance workload, mostly via normal births
 - ▶ Challenge: Workload is Operational, Cognitive, Emotional
 - ▶ Operational: Work content of a task, in time-units
 - ▶ Emotional: e.g. Mother and fetus-in-stress, suddenly fetus dies
- ⇒ Need help: A. Rafaeli & students (Psychology) - Ongoing

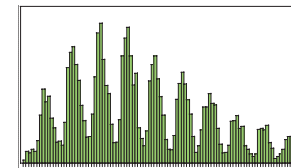
41

LogNormal & Beyond: Length-of-Stay in a Hospital

Israeli Hospital, in Days: LN



Israeli Hospital, in Hours: Mixture



Explanation: Patients released around 3pm (1pm in Singapore)

Why Bother ?

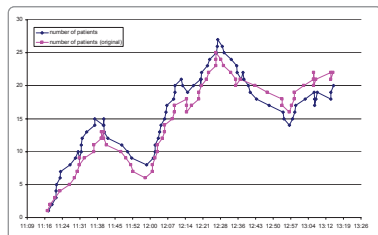
- ▶ Hourly Scale: Staffing, ...
- ▶ Daily: Flow / Bed Control, ...

42

Prerequisite II: Models (Fluid Q's)

"Laws of Large Numbers" capture Predictable Variability
Deterministic Models: Scale Averages-out Stochastic Individualism

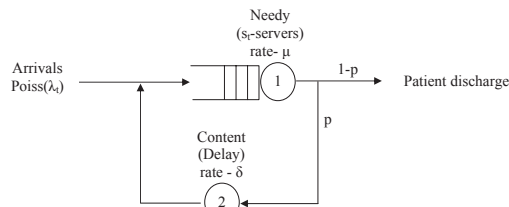
Severely-Wounded Patients, 11:00-13:00 (Censored LOS)



- ▶ Paths of doctors, nurses, patients (100+, 1 sec. resolution)
eg. (could) Help predict "What if 150+ casualties severely wounded ?"
- ▶ Transient Q's:
 - ▶ Control of Mass Casualty Events (w/ I. Cohen, N. Zychlinski)
 - ▶ Chemical MCE = Needy-Content Cycles (w/ G. Yom-Tov)

43

The Basic Service-Network Model: Erlang-R



Erlang-R (IE: Repairman Problem 50's; CS: Central-Server 60's) =
2-station "Jackson" Network = (M/M/S, M/M/∞) :

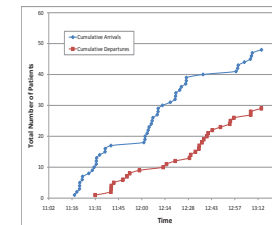
- ▶ $\lambda(t)$ – Time-Varying Arrival rate
- ▶ $S(\cdot)$ – Number of Servers (Nurses / Physicians).
- ▶ μ – Service rate ($E[\text{Service}] = \frac{1}{\mu}$)
- ▶ p – ReEntrant (Feedback) fraction
- ▶ δ – Content-to-Needy rate ($E[\text{Content}] = \frac{1}{\delta}$)

44

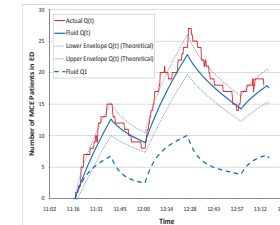
Erlang-R: Fitting a Simple Model to a Complex Reality

Chemical MCE Drill (Israel, May 2010)

Arrivals & Departures (RFID)



Erlang-R (Fluid, Diffusion)



- ▶ Recurrent/Repeated services in MCE Events: eg. Injection every 15 minutes
- ▶ Fluid (Sample-path) Modeling, via Functional Strong Laws of Large Numbers
- ▶ Stochastic Modeling, via Functional Central Limit Theorems
 - ▶ ED in MCE: Confidence-interval, usefully narrow for Control
 - ▶ ED in normal (time-varying) conditions: Personnel Staffing

45

Prerequisite II: Models (Diffusion/QED's Q's)

Traditional Queueing Theory predicts that **Service-Quality** and **Servers' Efficiency** must be traded off against each other.

For example, **M/M/1** (single-server queue): **91%** server's utilization goes with

$$\text{Congestion Index} = \frac{E[\text{Wait}]}{E[\text{Service}]} = 10.$$

and only **9%** of the customers are served immediately upon arrival.

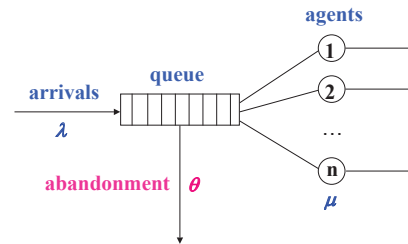
Yet, **heavily-loaded** queueing systems with **Congestion Index = 0.1** (Waiting one order of magnitude less than Service) are prevalent:

- ▶ **Call Centers:** Wait "**seconds**" for **minutes** service;
- ▶ **Transportation:** Search "**minutes**" for **hours** parking;
- ▶ **Hospitals:** Wait "**hours**" in ED for **days** hospitalization in IW's;

and, moreover, a significant fraction are not delayed in queue. (For example, in well-run call-centers, **50%** served "immediately", along with over **90%** agents' utilization, is not uncommon) ? **QED**

46

The Basic Staffing Model: Erlang-A (M/M/N + M)



Erlang-A (Palm 1940's) = **Birth & Death Q**, with parameters:

- ▶ λ – **Arrival** rate (Poisson)
- ▶ μ – **Service** rate (Exponential; $E[S] = \frac{1}{\mu}$)
- ▶ θ – **Patience** rate (Exponential, $E[\text{Patience}] = \frac{1}{\theta}$)
- ▶ n – Number of **Servers** (Agents).

47

Testing the Erlang-A Primitives

- ▶ **Arrivals:** Poisson?
- ▶ **Service-durations:** Exponential?
- ▶ **(Im)Patience:** Exponential?
- ▶ Primitives independent (eg. Impatience and Service-Durations)?
- ▶ Customers / Servers Homogeneous?
- ▶ Service discipline FCFS?
- ▶ ... ?

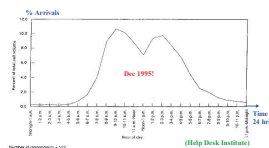
Validation: Support? Refute?

48

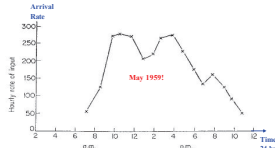
Arrivals to Service

Arrival-Rates to Three Call Centers

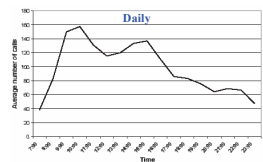
Dec. 1995 (U.S. 700 Helpdesks)



May 1959 (England)



November 1999 (Israel)



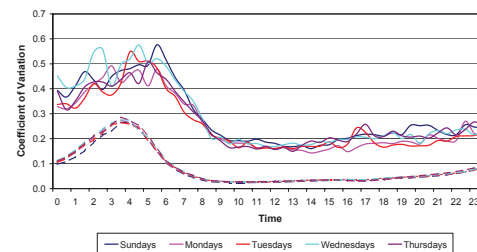
Random Arrivals "must be"
(Axiomatically)
Time-Inhomogeneous Poisson

49

Arrivals to Service: only Poisson-Relatives

Arrival-Counts: **Coefficient-of-Variation (CV)**, per 30 min.

Israeli-Bank Call-Center, 263 regular days (4/2007 - 3/2008)

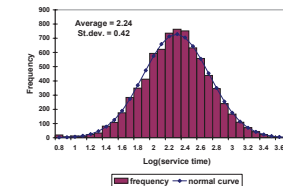


- ▶ **Poisson CV** (Dashed Line) = $1/\sqrt{\text{mean arrival-rate}}$
- ▶ Poisson CV's \ll **Sampled CV's** (Solid) \Rightarrow **Over-Dispersion**
- \Rightarrow **Modeling** (Poisson-Mixture) of and **Staffing** ($> \sqrt{\cdot}$) against **Time-Varying Over-Dispersed Arrivals** (w/ S. Maman & S. Zeltyn)

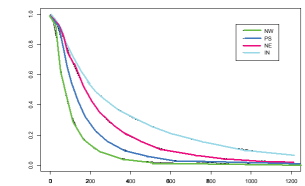
50

Service Durations: LogNormal Prevalent

Israeli Bank Log-Histogram



Service-Classes Survival-Functions

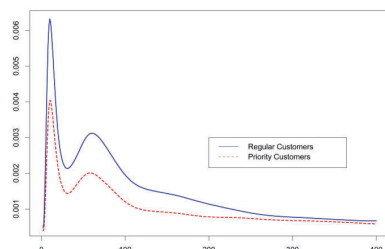


- **New Customers:** **2 min** (NW);
- **Regulars:** **3 min** (PS);
- **Stock:** **4.5 min** (NE);
- **Tech-Support:** **6.5 min** (IN).
- ▶ Service Durations are **LogNormal (LN)** and **Heterogeneous**

51

(Im)Patience while Waiting (Palm 1943-53)

Hazard Rate of (Im)Patience Distribution \propto Irritation
Regular over VIP Customers – Israeli Bank

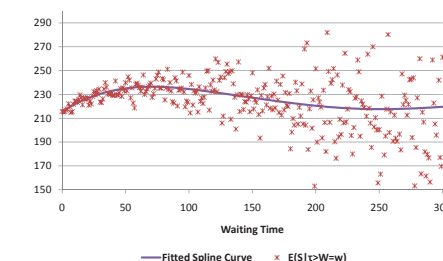


- ▶ **VIP** Customers are **more Patient** (Needy)
- ▶ **Peaks** of abandonment at times of **Announcements**
- ▶ Challenges: **Un-Censoring, Dependence (vs. KM), Smoothing**
- requires **Call-by-Call Data**

52

Dependent Primitives: Service- vs. Waiting-Time

Average Service-Time as a function of Waiting-Time
U.S. Bank, Retail, Weedays, January-June, 2006



\Rightarrow Focus on (**Patience, Service-Time**) jointly, w/ Reich and Ritov.
 $E[S | \text{Patience} = w]$, $w \geq 0$: **Service-Time of the Unserved.**

53

Erlang-A: Practical Relevance?

Experience:

- ▶ Arrival process **not pure Poisson** (time-varying, σ^2 too large)
- ▶ Service times **not Exponential** (typically close to LogNormal)
- ▶ Patience times **not Exponential** (various patterns observed).
- ▶ Building Blocks need **not be independent** (eg. long wait associated with long service; with w/ **M. Reich** and **Y. Ritov**)
- ▶ Customers and Servers **not homogeneous** (classes, skills)
- ▶ Customers return for service (after busy, abandonment; dependently; **P. Khudiakov**, **M. Gorfine**, **P. Feigin**)
- ▶ ... , and more.

Question: **Is Erlang-A Relevant?**

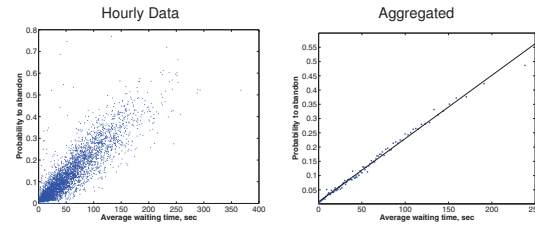
YES ! Fitting a Simple Model to a Complex Reality, both **Theoretically** and **Practically**

54

Estimating (Im)Patience: via $P\{Ab\} \propto E[W_q]$

"Assume" **Exp**(θ) (im)patience. Then, $P\{Ab\} = \theta \cdot E[W_q]$.

% Abandonment vs. Average Waiting-Time Bank Anonymous (JASA): Yearly Data



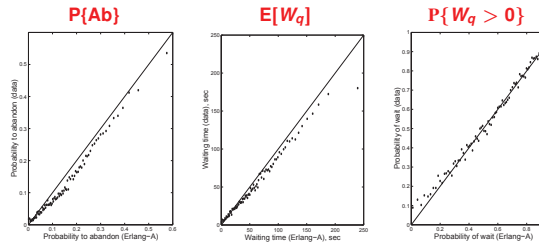
Graphs based on 4158 hour intervals.

Estimate of mean (im)patience: 250/0.55 sec. \approx **7.5 minutes**.

55

Erlang-A: Fitting a Simple Model to a Complex Reality

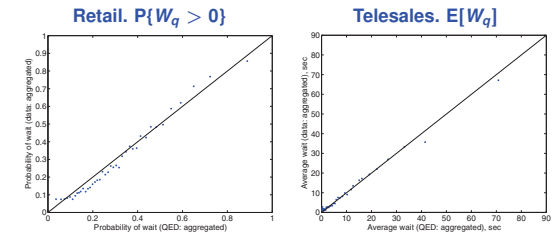
- Bank Anonymous Small Israeli Call-Center
- (Im)Patience (θ) estimated via $P\{Ab\} / E[W_q]$
- Graphs: **Hourly Performance vs. Erlang-A Predictions**, during 1 year (aggregating groups with 40 similar hours).



56

Erlang-A: Fitting a Simple Model to a Complex Reality

Large U.S. Bank



Partial success – in **some** cases Erlang-A **does not work** well (Networking, SBR).

Ongoing **Validation** Project, w/ Y. Nardi, O. Plonsky, S. Zeltyn

57

Erlang-A: Simple, but Not Too Simple

Practical (Data-Based) questions, started in **Brown et al. (JASA)**:

- Fitting Erlang-A (**Validation**, w/ Nardi, Plonsky, Zeltyn).
- Why does it practically work? justify **robustness**.
- When does it fail? chart **boundaries**.
- Generate needs for **new theory**.

Theoretical Framework: Asymptotic Analysis, as load- and staffing-levels increase, which reveals model-essentials:

- Efficiency-Driven (ED)** regime: Fluid models (deterministic)
- Quality- and Efficiency-Driven (QED)**: Diffusion refinements.

Motivation: Moderate-to-large service systems (100's - 1000's servers), notably **Call-Centers**.

Results turn out **accurate** enough to also cover <10 servers:

- Practically Important:** Relevant to **Healthcare** (First: F. de Véricourt and O. Jennings; w/ G. Yom-Tov; Y. Marmor, S. Zeltyn; H. Kaspi, I. Zaeid)
- Theoretically Justifiable:** Gap-Analysis by A. Janssen, J. van Leeuwen, B. Zhang, B. Zwart.

58

Operational Regimes: Conceptual Framework

R: Offered Load

Def. R = Arrival-rate \times Average-Service-Time = $\frac{\lambda}{\mu}$

eg. R = 25 calls/min. \times 4 min./call = **100**

N = #Agents ? **Intuition**, as R or N increase unilaterally.

QD Regime: $N \approx R + \delta R$, $0.1 < \delta < 0.25$ (eg. $N = 115$)

- Framework developed in **O. Garnett's** MSc thesis
- Rigorously: $(N - R)/R \rightarrow \delta$, as $N, \lambda \uparrow \infty$, with μ fixed.
- Performance: Delays are rare events

ED Regime: $N \approx R - \gamma R$, $0.1 < \gamma < 0.25$ (eg. $N = 90$)

- Essentially **all** customers are delayed
- Wait same order as service-time; $\gamma\%$ Abandon (10-25%).

QED Regime: $N \approx R + \beta\sqrt{R}$, $-1 < \beta < +1$ (eg. $N = 100$)

- Erlang 1913-24, **Halfin & Whitt** 1981 (for Erlang-C)
- %Delayed between 25% and 75%
- $E[\text{Wait}] \propto \frac{1}{\sqrt{N}} \times E[\text{Service}]$ (**sec vs. min**); 1-5% Abandon.

59

Operational Regimes: Rules-of-Thumb, w/ S. Zeltyn

Constraint	P{Ab}		E[W]		P{W > T}	
	Tight	Loose	Tight	Loose	Tight	Loose
	1-10%	$\geq 10\%$	$\leq 10\%E[\tau]$	$\geq 10\%E[\tau]$	$0 \leq T \leq 10\%E[\tau]$	$T \geq 10\%E[\tau]$
Offered Load					$5\% \leq \alpha \leq 50\%$	$5\% \leq \alpha \leq 50\%$
Small (10's)	QED	QED	QED	QED	QED	QED
Moderate-to-Large (100's-1000's)	QED	ED, QED	ED, QED if $\tau \not\leq \exp$	QED	ED+QED	ED+QED

ED: $N \approx R - \gamma R$ ($0.1 \leq \gamma \leq 0.25$).

QD: $N \approx R + \delta R$ ($0.1 \leq \delta \leq 0.25$).

QED: $N \approx R + \beta\sqrt{R}$ ($-1 \leq \beta \leq 1$).

ED+QED: $N \approx (1 - \gamma)R + \beta\sqrt{R}$ (γ, β as above).

WFM: How to determine specific staffing level N ? e.g. β .

60

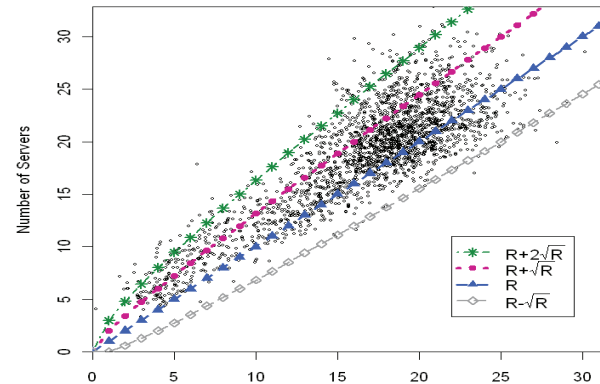
Operational Regimes: Scaling, Performance, w/ I. Gurvich & J. Huang

Erlang-A	Conventional scaling			MS scaling			NDS scaling		
	Sub	Critical	Super	QD	QED	ED	ED+QED	Sub	Super
Offered load per server	$\frac{R}{n} < 1$	$1 - \frac{R}{n} \approx 1$	$\frac{R}{n} > 1$	$\frac{R}{n}$	$1 - \frac{R}{n}$	$\frac{R}{n}$	$\frac{R}{n} - \beta\sqrt{\frac{R}{n}}$	$\frac{R}{n}$	$1 - \frac{R}{n}$
Arrival rate λ	$\frac{\lambda}{n}$	$\mu - \frac{\lambda}{n}\mu$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$	$\mu - \beta\sqrt{\frac{\lambda}{n}}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n} - \beta\sqrt{\frac{\lambda}{n}}$	$\frac{\lambda}{n}$	$\frac{\lambda}{n}$
Number of servers		1			n				n
Time-scale		n			1				n
Abandonment rate		θ/n			θ				θ/n
Staffing level	$\frac{2}{2(1+\delta)}$	$\frac{2}{2(1+\delta)}$	$\frac{2}{2(1-\gamma)}$	$\frac{2}{2(1+\delta)}$	$\frac{2}{2(1-\gamma)}$	$\frac{2}{2(1-\gamma)}$	$\frac{2}{2(1-\gamma)}$	$\frac{2}{2(1+\delta)}$	$\frac{2}{2(1-\gamma)}$
Utilization	$\frac{1}{1+\delta}$	$1 - \frac{\sqrt{2\lambda\theta}}{\sqrt{2\lambda\theta}}$	1	$\frac{1}{1+\delta}$	$1 - \frac{\sqrt{2\lambda\theta}}{\sqrt{2\lambda\theta}}$	1	$1 - \frac{\sqrt{2\lambda\theta}}{\sqrt{2\lambda\theta}}$	$\frac{1}{1+\delta}$	1
$E[Q]$	$\frac{1}{\mu}$	$\frac{1}{\mu} \sqrt{\frac{2\lambda\theta}{\mu}}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu} \sqrt{\frac{2\lambda\theta}{\mu}}$	$\frac{1}{\mu}$	$\frac{1}{\mu} \sqrt{\frac{2\lambda\theta}{\mu}}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$
$E[A]$	$\frac{1}{\mu}$	$\frac{1}{\mu} \sqrt{\frac{2\lambda\theta}{\mu}}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu} \sqrt{\frac{2\lambda\theta}{\mu}}$	$\frac{1}{\mu}$	$\frac{1}{\mu} \sqrt{\frac{2\lambda\theta}{\mu}}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$
$P\{W_q > 0\}$	$\alpha_0 \in (0,1)$	≈ 1	≈ 1	$\alpha_0 \in (0,1)$	≈ 1	≈ 1	≈ 1	$\alpha_0 \in (0,1)$	≈ 1
$P\{W_q > T\}$	$\alpha_0 e^{-\delta\theta T}$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$	≈ 0	$\frac{G(T)}{G(0)}$	α_0 if $G(T) = \gamma$	≈ 0	$\frac{G(T)}{G(0)}$	$1 + O(\frac{1}{n})$
Congestion	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$	$\frac{1}{\mu}$

61

QED Call Center: Staffing (N) vs. Offered-Load (R)

IL Telecom; June-September, 2004; w/ Nardi, Plonski, Zeltyn

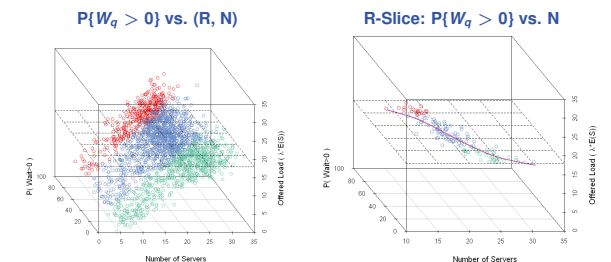


2205 half-hour intervals in an Israeli Call Center

62

QED Call Center: Performance

Large Israeli Bank



3 Operational Regimes:

- QD:** $\leq 25\%$
- QED:** 25% – 75%
- ED:** $\geq 75\%$

63

QED Theory (Erlang '13; Halfin-Whitt '81; Garnett MSc; Zeltyn PhD)

Consider a sequence of **steady-state** M/M/ N + G queues, $N = 1, 2, 3, \dots$. Then the following points of view are **equivalent**, as $N \uparrow \infty$:

- **QED** $\% \{\text{Wait} > 0\} \approx \alpha$, $0 < \alpha < 1$;
- **Customers** $\% \{\text{Abandon}\} \approx \frac{\gamma}{\sqrt{N}}$, $0 < \gamma$;
- **Agents** $\text{OCC} \approx 1 - \frac{\beta + \gamma}{\sqrt{N}}$, $-\infty < \beta < \infty$;
- **Managers** $N \approx R + \beta \sqrt{R}$, $R = \lambda \times E(S)$ not small;

- **QED performance: Laplace Method** (asymptotics of integrals).
- **Parameters:** Arrivals and Staffing - β , Services - μ , (Im)Patience - $g(0)$ = **patience density at the origin**.

64

Erlang-A: QED Approximations (Examples)

Assume **Offered Load** R not small ($\lambda \rightarrow \infty$).

Let $\hat{\beta} = \beta \sqrt{\frac{\mu}{\theta}}$, $h(\cdot) = \frac{\phi(\cdot)}{1 - \Phi(\cdot)}$ = hazard rate of $\mathcal{N}(0, 1)$.

- **Delay Probability:**

$$P\{W_q > 0\} \approx \left[1 + \sqrt{\frac{\theta}{\mu}} \cdot \frac{h(\hat{\beta})}{h(-\hat{\beta})} \right]^{-1}.$$

- **Probability to Abandon:**

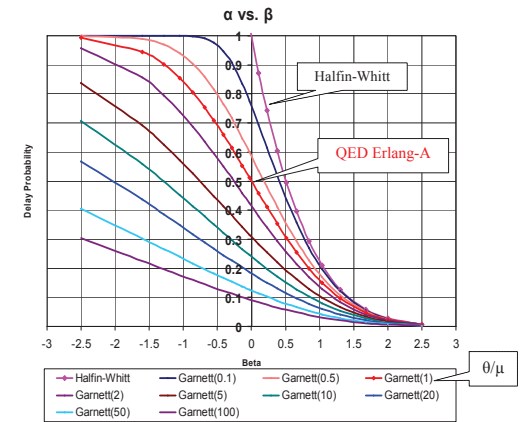
$$P\{\text{Ab}|W_q > 0\} \approx \frac{1}{\sqrt{n}} \cdot \sqrt{\frac{\theta}{\mu}} \cdot [h(\hat{\beta}) - \hat{\beta}].$$

- $P\{\text{Ab}\} \propto E[W_q]$, both order $\frac{1}{\sqrt{n}}$:

$$\frac{P(\text{Ab})}{E[W_q]} = \theta.$$

65

Garnett / Halfin-Whitt Functions: $P\{W_q > 0\}$



66

QED Intuition: Why $P\{W_q > 0\} \in (0, 1)$?

- Why **subtle**: Consider a large service system (e.g. call center).
 - Fix λ and let $n \uparrow \infty$: $P\{W_q > 0\} \downarrow 0$.
 - Fix n and let $\lambda \uparrow \infty$: $P\{W_q > 0\} \uparrow 1$.
 - ⇒ **Must** have **both** λ and n increase simultaneously:
 - ⇒ (CLT) **Square-root staffing**: $n \approx R + \beta \sqrt{R}$.

- Erlang-A** (M/M/n+M), with parameters λ, μ, θ, n , in which $\mu = \theta$: (Im)Patience and Service-times are equally distributed.

- Steady-state: $L(M/M/n+M) \stackrel{d}{=} L(M/M/\infty) \stackrel{d}{=} \text{Poisson}(R)$, with $R = \lambda/\mu$ (Offered-Load)
- $\text{Poisson}(R) \stackrel{d}{\approx} R + Z\sqrt{R}$, with $Z \stackrel{d}{=} \mathcal{N}(0, 1)$.
- $P\{W_q(M/M/n+M) > 0\} \stackrel{\text{PASTA}}{=} P\{L(M/M/n+M) \geq n\} \stackrel{\mu=\theta}{=} P\{L(M/M/\infty) \geq n\} \approx P\{R + Z\sqrt{R} \geq n\} =$

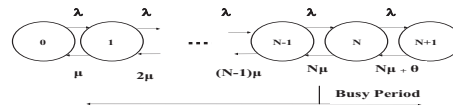
$$P\{L(M/M/\infty) \geq n\} \approx P\{R + Z\sqrt{R} \geq n\} =$$

$$P\{Z \geq (n - R)/\sqrt{R}\} \stackrel{\sqrt{\cdot} \text{ staffing}}{\approx} P\{Z \geq \beta\} = 1 - \Phi(\beta).$$

- QED **Excursions**

67

QED Intuition via Excursions: Busy-Idle Cycles



$Q(0) = N$: all servers busy, no queue.

Let $T_{N,N-1} = E[\text{Busy Period}]$ down-crossing $N \downarrow N-1$

$T_{N-1,N} = E[\text{Idle Period}]$ up-crossing $N-1 \uparrow N$

Then $P(\text{Wait} > 0) = \frac{T_{N,N-1}}{T_{N,N-1} + T_{N-1,N}} = \left[1 + \frac{T_{N-1,N}}{T_{N,N-1}} \right]^{-1}.$

68

QED Intuition via Excursions: Asymptotics

Calculate $T_{N-1,N} = \frac{1}{\lambda_N E_{1,N-1}} \sim \frac{1}{N\mu \times h(-\beta)/\sqrt{N}} \sim \frac{1}{\sqrt{N}} \cdot \frac{1/\mu}{h(-\beta)}$

$T_{N,N-1} = \frac{1}{N\mu\pi + (0)} \sim \frac{1}{\sqrt{N}} \cdot \frac{\beta/\mu}{h(\delta)/\delta}$, $\delta = \beta\sqrt{\mu}/\theta$

Both apply as $\sqrt{N}(1 - \rho_N) \rightarrow \beta$, $-\infty < \beta < \infty$.

Hence, $P(\text{Wait} > 0) \sim \left[1 + \frac{h(\delta)/\delta}{h(-\beta)/\beta} \right]^{-1}.$

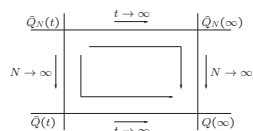
69

Process Limits (Queueing, Waiting)

- $\tilde{Q}_N = \{\tilde{Q}_N(t), t \geq 0\}$: **stochastic process** obtained by centering and rescaling:

$$\tilde{Q}_N = \frac{Q_N - N}{\sqrt{N}}$$

- $\tilde{Q}_N(\infty)$: stationary distribution of \tilde{Q}_N
- $\tilde{Q} = \{\tilde{Q}(t), t \geq 0\}$: process defined by: $\tilde{Q}_N(t) \xrightarrow{d} \tilde{Q}(t)$.



Approximating (Virtual) **Waiting Time**

$$\tilde{V}_N = \sqrt{N} V_N \Rightarrow \tilde{V} = \left[\frac{1}{\mu} \tilde{Q} \right]^+_{\tilde{V}_0}$$

QED Erlang-X (Markovian Q's: Performance Analysis)

- Pre-History, 1914: **Erlang** (Erlang-B = M/M/n/n, Erlang-C = M/M/n)
- Pre-History, 1974: Jagerman (Erlang-B)
- History Milestone, 1981: **Halfin-Whitt** (Erlang-C, GI/M/n)
- Erlang-A (M/M/N+M), 2002: w/ **Garnett** & Reiman
- Erlang-A with General (Im)Patience (M/M/N+G), 2005: w/ Zeltyn
- Erlang-C (ED+QED), 2009: w/ Zeltyn
- Erlang-B with Retrial, 2010: Avram, Janssen, van Leeuwen
- Refined Asymptotics (Erlang A/B/C), 2008-2011: Janssen, van Leeuwen, Zhang, Zwart
- NDS Erlang-C/A, 2009: Atar
- Production Q's, 2011: Reed & Zhang
- Universal Erlang-R, ongoing: w/ Gurvich & Huang
- Queueing Networks:
 - (Semi-)Closed: Nurse Staffing (Jennings & de Vericourt), CCs with IVR (w/ Khudiyakov), Erlang-R (w/ Yom-Tov)
 - CCs with Abandonment and Retrials: w. Massey, Reiman, Rider, Stolyar
 - Markovian Service Networks: w/ Massey & Reiman
- Leaving out:
 - **Non-Exponential Service Times**: M/D/n (Erlang-D), G/Ph/n, ..., G/GI/n+GI, Measure-Valued Diffusions
 - **Dimensioning** (Staffing): M/M/n, ..., time-varying Q's, V- and Reversed-V, ...
 - **Control**: V-network, Reversed-V, ..., SBRNets

71

Back to "Why does Erlang-A Work?"

Theoretical (Partial) Answer:

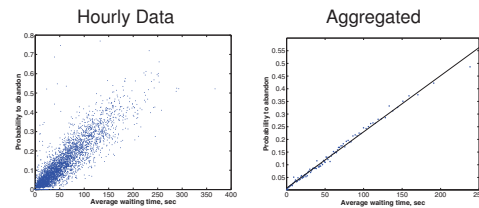
$$M_t^{i,j} / G^* / N_t + G \stackrel{d}{\approx} (M/M/N+M)_t, \quad t \geq 0.$$

- **Over-Dispersed Arrivals**: $R + \beta R^c$, c-Staffing ($c \geq 1/2$).
- **General Patience**: Behavior at the origin matters most (only).
- **General Services**: Empirical insensitivity beyond the mean.
- **Heterogeneous Customers / Servers**: State-Collapse.
- **Time-Varying Arrivals**: Modified Offered-Load approximations.
- **Dependent Building-Blocks**: eg. When (Im)Patience and Service-Times correlated (positively):
 - Predict performance with $E[S | \text{Served}]$.
 - Calculate offered-load with $E[S] = E[S | \text{Wait} = 0]$.
 - Note: staffing \leftarrow service-times \leftarrow waiting / abandonment \leftarrow staffing

72

“Why does Erlang-A Work?” General Patience

Israeli Bank: Yearly Data



Theory:

Erlang-A: $P\{Ab\} = \theta \cdot E[W_q]$;

M/M/N+G: $P\{Ab\} \approx g(0) \cdot E[W_q]$.

$g(0)$ = Patience-density at origin

Recipe:

In both cases, use Erlang-A, with $\hat{\theta} = \widehat{P\{Ab\}} / \widehat{E[W_q]}$ (slope above).

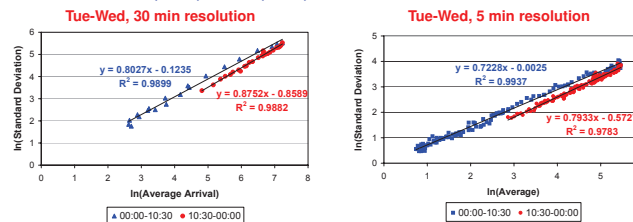
References on $g(0)$:

- Stationary M/M/N+GI, w/ Zeltyn
- Process G/GI/N+GI: w/ Momcilovic; Dai & He;

73

“Why does Erlang-A Work?” Over-Dispersion

$\ln(STD)$ vs. $\ln(AVG)$ (Israeli Bank, 4/2007-3/2008)



Significant linear relations (w/ Aldor & Feigin; then w/

Maman & Zeltyn):

$$\ln(STD) = c \cdot \ln(AVG) + a$$

(Poisson: $STD = AVG^{1/2}$, hence $c = 1/2, a = 0$.)

74

Over-Dispersion: Random Arrival-Rates

Linear relation between $\ln(STD)$ and $\ln(AVG)$ gives rise to:

Poisson-Mixture (Doubly-Poisson, Cox) model for Arrivals:
Poisson(Λ) with **Random-Rate** of the form

$$\Lambda = \lambda + \lambda^c \cdot X, \quad c \leq 1;$$

- ▶ c determines magnitude of over-dispersion (λ^c)
 $c = 1$, proportional to λ ; $c \leq 1/2$, Poisson-level;
- In **Call Centers**: $c \approx 0.75 - 0.85$ (significant over-dispersion).
- In **Emergency Departments**, $c \approx 0.5$ (Poisson).
- ▶ X random-variable with $E[X] = 0$ ($E[\Lambda] = \lambda$), capturing the magnitude of **stochastic deviation** from mean arrival-rate: under conventional Gamma prior (λ large), X can be taken Normal with std. derived from the intercept.

QED-c Regime: Erlang-A, with Poisson(Λ) arrivals, amenable to asymptotic analysis (with S. Maman & S. Zeltyn)

75

Over-Dispersion: The QED-c Regime

QED-c Staffing: Under offered-load $R = \lambda \cdot E[S]$,

$$N = R + \beta \cdot R^c, \quad 0.5 < c < 1$$

Performance measures (M/M/N + G):

- Delay probability: $P\{W_q > 0\} \sim 1 - G(\beta)$
- Abandonment probability: $P\{Ab\} \sim \frac{E[X - \beta]_+}{n^{1-c}}$
- Average offered wait: $E[V] \sim \frac{E[X - \beta]_+}{n^{1-c} \cdot g_0}$
- Average actual wait: $E_{\Lambda, N}[W] \sim E_{\Lambda, N}[V]$

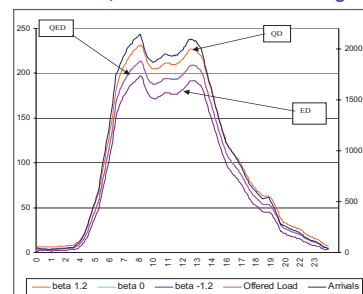
76

Why Does Erlang-A Work? Time-Varying Arrival Rates

Square-Root Staffing: $N_t = R_t + \beta \sqrt{R_t}$, $-\infty < \beta < \infty$

What is R_t , the **Offered-Load** at time t ? ($R_t \neq \lambda_t \times E[S]$)

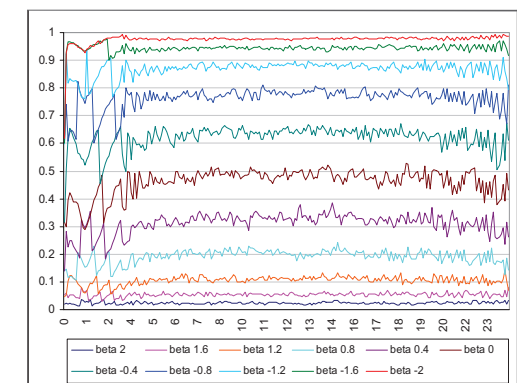
Arrivals, Offered-Load and Staffing



77

Time-Stable Performance of Time-Varying Systems

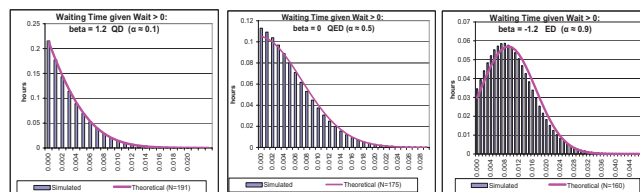
Delay Probability = As in the **Stationary Erlang-A** (Garnett)



78

Time-Stable Performance of Time-Varying Systems

Waiting Time, Given Waiting: Empirical vs. Theoretical Distribution



- **Empirical:** Simulate **time-varying** $M_t/M/N_t + M$ ($\lambda_t, N_t = R_t + \beta \sqrt{R_t}$)
- **Theoretical:** Naturally-corresponding **stationary** Erlang-A, with QED β -staffing (some **Averaging** Principle?)
- **Generalizes** up to a single-station within a complex network (eg. Doctors in an Emergency Department).

79

What is the Offered-Load $R(t)$?

- ▶ Offered-Load Process: $L(\cdot)$ = **Least** number of **servers** that guarantees **no delay**.
- ▶ **Offered-Load** Function $R(t) = E[L(t)]$, $t \geq 0$.
Think $M_t/G/N_t^2 + G$ vs. $M_t/G/\infty$: **Ample-Servers**.

Four (all useful) representations, capturing “workload before t ”:

$$R(t) = E[L(t)] = \int_{-\infty}^t \lambda(u) \cdot P(S > t - u) du = E \left[A(t) - A(t - S) \right] = E \left[\int_{t-S}^t \lambda(u) du \right] = E[\lambda(t - S_e)] \cdot E[S] \approx \dots$$

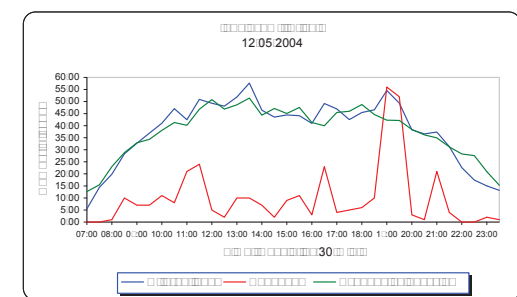
- ▶ $\{A(t), t \geq 0\}$ Arrival-Process, rate $\lambda(\cdot)$;
- ▶ $S(S_e)$ generic Service-Time (Residual Service-Time).
- ▶ Relating L, λ, S (“ W ”): **Time-Varying Little’s Formula**.
Stationary models: $\lambda(t) \equiv \lambda$ then $R(t) \equiv \lambda \times E[S]$.

QED-c: $N_t = R_t + \beta R_t^c$, $1/2 \leq c < 1$; ($c = 1$ separate analysis).

80

The Offered-Load $R(t)$, $t \geq 0$

- ▶ **Backbone** of time-varying staffing:
 - ▶ Practically **robust**: up to a station within a complex network (ED).
 - ▶ Theoretically **challenging**: only Erlang-A with $\mu = \theta$ tractable.
- ▶ Process: $L(\cdot)$ = **Least** number of **servers** that guarantees **no delay**.
- ▶ **Offered-Load** Function $R(\cdot) = E[L(\cdot)]$ ($M_t/G/N_t^2 + G \leftrightarrow M_t/G/\infty$).



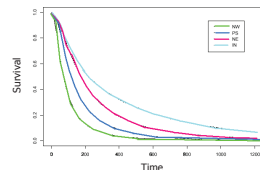
81

Estimating / Predicting the Offered-Load

Must account for “**service times of abandoning customers**”.

- ▶ Prevalent Assumption: Services and (Im)Patience independent.
- ▶ But recall Patient VIPs: Willing to wait more for longer services.

Survival Functions by Type (Small Israeli Bank)



Service times stochastic order: $S_{New} \stackrel{st}{<} S_{Reg} \stackrel{st}{<} S_{VIP}$

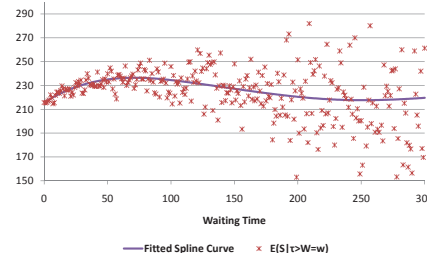
Patience times stochastic order: $T_{New} \stackrel{st}{<} T_{Reg} \stackrel{st}{<} T_{VIP}$

82

Dependent Primitives: Service- vs. Waiting-Time

Average Service-Time as a function of Waiting-Time

U.S. Bank, Retail, Weekdays, January-June, 2006



⇒ Focus on (**Patience, Service-Time**) jointly , w/ Reich and Ritov.
 $E[S | \text{Patience} = w], w \geq 0$: **Service-Time of the Unserved.**

83

(Imputing) Service-Times of Abandoning Customers

w/ M. Reich, Y. Ritov:

1. **Estimate** $g(w) = E[S | \text{Patience} > \text{Wait} = w], w \geq 0$:
Mean service time of those **served after waiting exactly** w units of time (via non-linear regression: $S_i = g(W_i) + \varepsilon_i$);

2. **Calculate**

$$E[S | \text{Patience} = w] = g(w) - \frac{g'(w)}{h_r(w)};$$

$h_r(w)$ = hazard-rate of (im)patience (via un-censoring);

3. **Offered-load** calculations: Impute $E[S | \text{Patience} = w]$ (or the conditional distribution).

Challenges: Stable and accurate inference of g, g', h_r .

84

Extending the Notion of the “Offered-Load”

- ▶ **Business** (Banking Call-Center): Offered **Revenues**
- ▶ **Healthcare** (Maternity Wards): Fetus in stress
 - ▶ 2 patients (Mother + Child) = high **operational** and **cognitive** load
 - ▶ Fetus dies ⇒ **emotional** load dominates
- ▶ ⇒
 - ▶ Offered **Operational** Load
 - ▶ Offered **Cognitive** Load
 - ▶ Offered **Emotional** Load
- ▶ ⇒ **Fair** Division of Load (Routing) between 2 Maternity Wards:
One attending to complications before birth, the other to complications after birth, and both share normal birth

85

The Technion SEE Center / Laboratory

Data-Based Service Science / Engineering



86