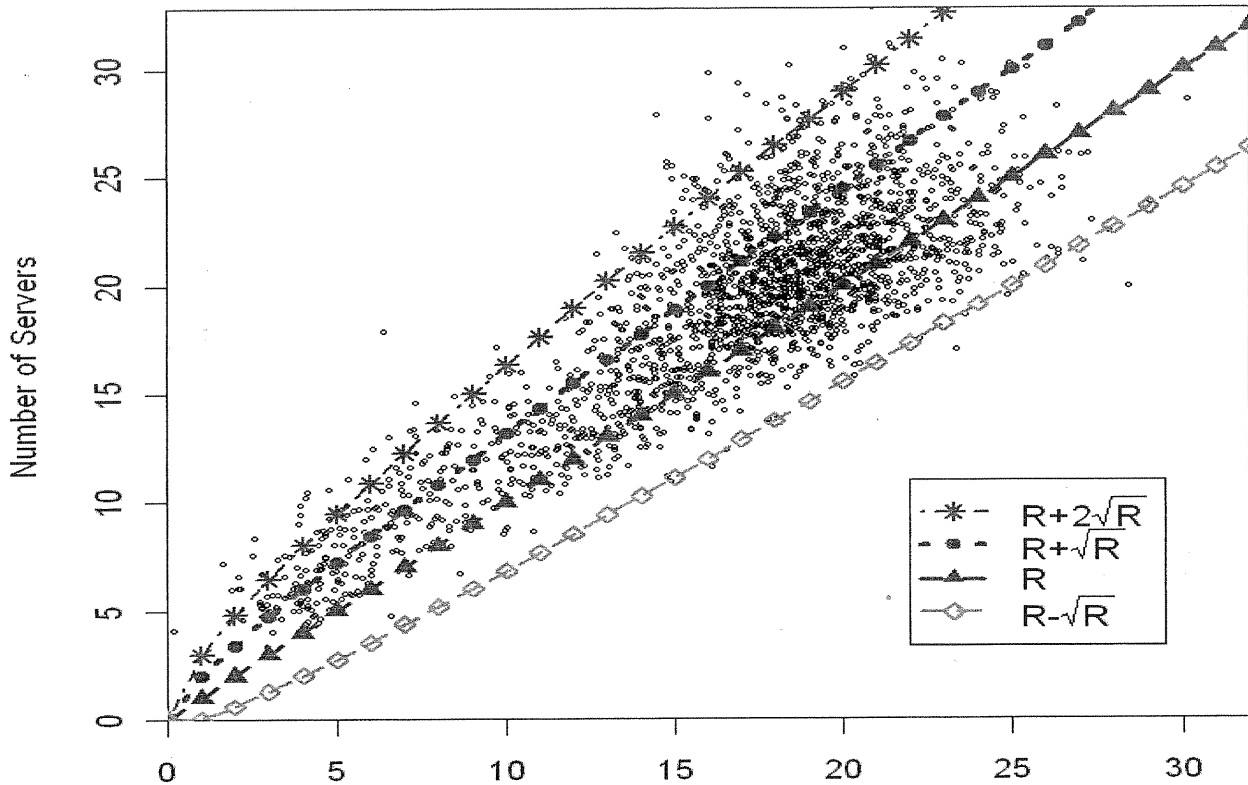


QED Call Center: Staffing (N) vs. Offered-Load (R)

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



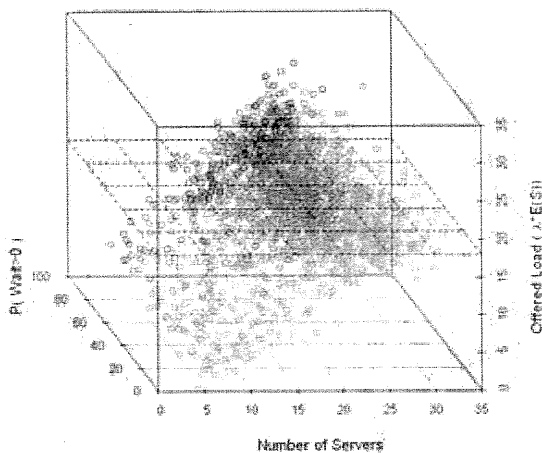
2205 half-hour intervals in an Israeli Call Center

7

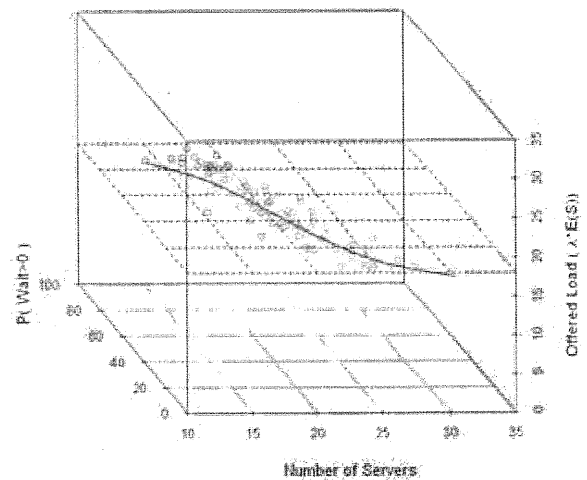
QED Call Center: Performance

Large Israeli Bank

$P\{W_q > 0\}$ vs. (R, N)



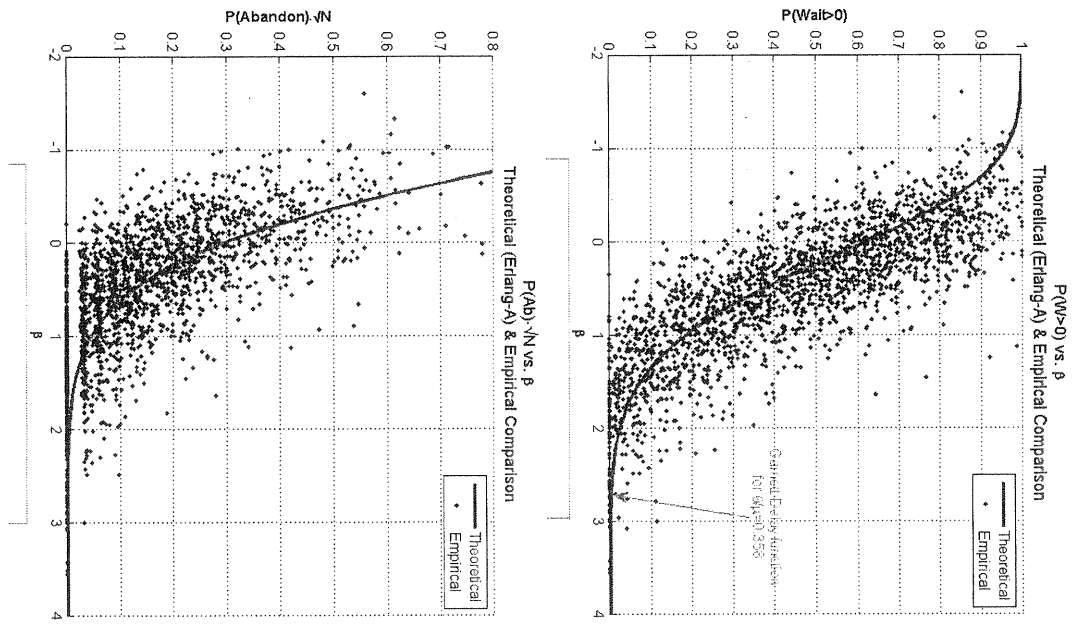
R-Slice: $P\{W_q > 0\}$ vs. N



3 Operational Regimes:

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

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



Erlang-A	μ fixed	Conventional scaling			NDS scaling			NDS scaling		
		Sub	Critical	Super	Sub	ED+QED	ED	Sub	Critical	Super
Offered load per server	$\frac{1}{1-\beta}$	$1 - \frac{\beta}{\mu}$	≈ 1	$\frac{1}{1-\beta} > 1$	$1 - \frac{\beta}{\mu}$	$1 - \frac{\beta}{\mu}$	$1 - \frac{\beta}{\mu}$	$1 - \frac{\beta}{\mu}$	$1 - \frac{\beta}{\mu}$	$1 - \frac{\beta}{\mu}$
Arrival rate λ	$\frac{1}{\mu}$	$\mu - \frac{\beta}{\mu}$	μ	$\mu - \frac{\beta}{\mu}$	$\mu - \frac{\beta}{\mu}$	$\mu - \frac{\beta}{\mu}$	$\mu - \frac{\beta}{\mu}$	$\mu - \frac{\beta}{\mu}$	$\mu - \frac{\beta}{\mu}$	$\mu - \frac{\beta}{\mu}$
Number of servers	1	1	1	1	1	1	1	1	1	1
Time-scale	n	n	n	n	n	n	n	n	n	n
Abandonment rate	θ/n	θ/n	θ/n	θ/n	θ/n	θ/n	θ/n	θ/n	θ/n	θ/n
Staffing level	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$	$\frac{n}{\lambda}(1+\delta)$
Utilization	$\frac{1}{1+\delta}$	$1 - \sqrt{\frac{\theta}{\mu h(\theta)}}$	1	1	1	1	1	1	1	1
E(Q)	$\frac{\theta}{\mu}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$
P(Ab)	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$	$\frac{n}{1+\delta}$
P(W _q > 0)	$\alpha_1 \in (0, 1)$	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1	≈ 1
P(W _q > T)	$\alpha_1 e^{-\frac{T}{1+\delta}}$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$	$1 + O(\frac{1}{n})$
Congestion $\frac{EW}{ES}$	$\alpha_1 \frac{\theta}{1+\delta}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$	$\sqrt{\frac{\theta}{\mu h(\theta)}}$

• $\delta > 0, \gamma \in (0, 1)$ and $\beta \in (-\infty, \infty)$;

• QD: $\theta = \frac{1}{1+\delta} e^{\frac{1}{1+\delta}}$ and $h(x) = \frac{\theta}{x}$;

• ED (ED+QED): $G(x^*) = \gamma$;

• QED: $\alpha_2 = 1 + \sqrt{\frac{\theta}{\mu h(\theta)}}$, here $\beta = \theta \sqrt{\frac{\theta}{\mu h(\theta)}}$ and $h(x) = \frac{\theta}{x}$;

• ED+QED: $\alpha_3 = G(T) \Phi(\sqrt{\frac{\theta}{\mu h(\theta)}}$;

• Conventional: critical: $P(W > T) = P(\frac{W}{\mu} > \frac{T}{\mu})$; super: $P(W > T) = P(\frac{W}{\mu} > \frac{T}{\mu})$; NDS: super: $P(W > T) = P(\frac{W}{\mu} > \frac{T}{\mu})$;

QED Erlang-A: Theoretical Motivation

QED staffing: $n \approx R + \beta\sqrt{R}$.

Assume $\theta = \mu$, namely “average service-time” = “average (im)patience”.

Recall and Note:

- If $\theta = \mu$, the number-in-system of $M/M/n+M$ has the same distribution of a corresponding $M/M/\infty$ (both are the same Birth&Death process). Formally, in steady-state:
 $L(M/M/n+M) \stackrel{d}{=} L(M/M/\infty)$.
- The steady-state distribution of $M/M/\infty$ with parameters λ and μ is **Poisson**(R), where $R = \lambda/\mu$ (offered-load).
- For R not too small, $\text{Poisson}(R)$ is approximately $\text{Normal}(R, R)$.
 Formally: $L(M/M/\infty) \stackrel{d}{\approx} R + Z\sqrt{R}$, where Z is standard normal.

We now use these facts to estimate the delay-probability for Erlang-A, in which $\theta = \mu$:

$$\begin{array}{ccc} P\{W_q(M/M/n+M) > 0\} & \stackrel{\text{PASTA}}{=} & P\{L(M/M/n+M) \geq n\} \\ & \stackrel{\theta=\mu}{=} & P\{L(M/M/\infty) \geq n\} \end{array}$$

Standardizing $L \approx R + Z\sqrt{R}$ reveals the QED regime, specifically how square-root staffing yields a non-degenerate delay-probability:

$$P\{W_q > 0\} \approx P\left\{Z \geq \frac{n-R}{\sqrt{R}}\right\} \approx 1 - \Phi(\beta).$$

The Erlang-A Queue in the QED-Regime

Theorem (with Garnett & Reiman, 2002)

The following **points of view** are equivalent:

- 0. QED:** $P\{W_q > 0\} \approx \alpha$, for some $0 < \alpha < 1$;
- 1. Manager:** $n \approx R + \beta\sqrt{R}$, for some $-\infty < \beta < \infty$;
- 2. Servers:** Occupancy $\approx 1 - \frac{\beta + \gamma}{\sqrt{n}}$;
- 3. Customers:** $P\{Ab\} \approx \frac{\gamma}{\sqrt{n}}$, for some $0 < \gamma < \infty$;

in which case

$$\alpha = \alpha\left(\beta, \frac{\mu}{\theta}\right) = \left[1 + \sqrt{\frac{\theta}{\mu}} \cdot \frac{h(\hat{\beta})}{h(-\beta)}\right]^{-1},$$

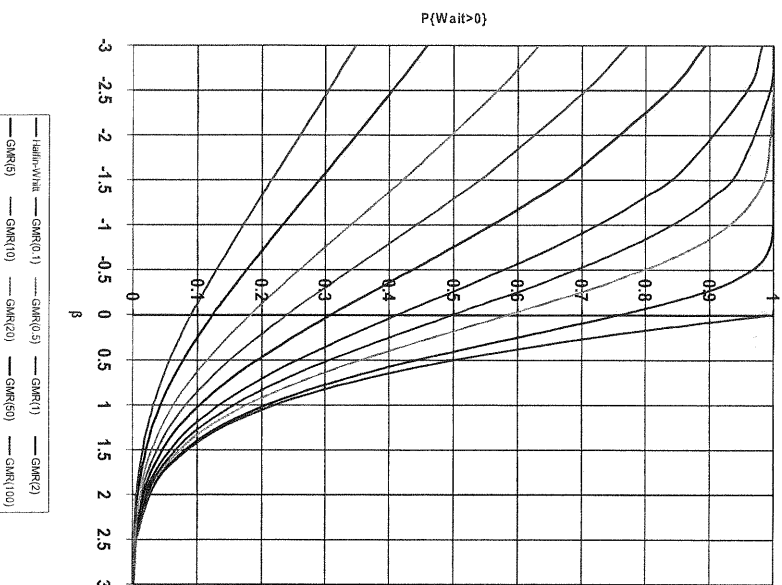
which we call the **Garnett Delay-Function**(s);

here $\hat{\beta} \triangleq \beta\sqrt{\frac{\mu}{\theta}}$, and

$$\gamma = \alpha \cdot \sqrt{\frac{\theta}{\mu}} \cdot [h(\hat{\beta}) - \hat{\beta}].$$

Erlang-A: The Garnett Delay-Functions

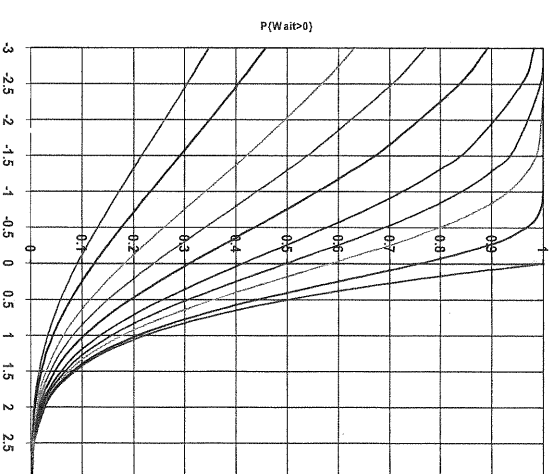
$P\{W_q > 0\}$ vs. the QOS parameter β , for varying patience θ/μ .



GMR(x) describes the asymptotic probability of delay as a function of β when $\theta/\mu = x$. Here, θ and μ are the abandonment and service rate, respectively.

Note: **Erlang-C** = limit of **Erlang-A**, as patience \uparrow indefinitely.

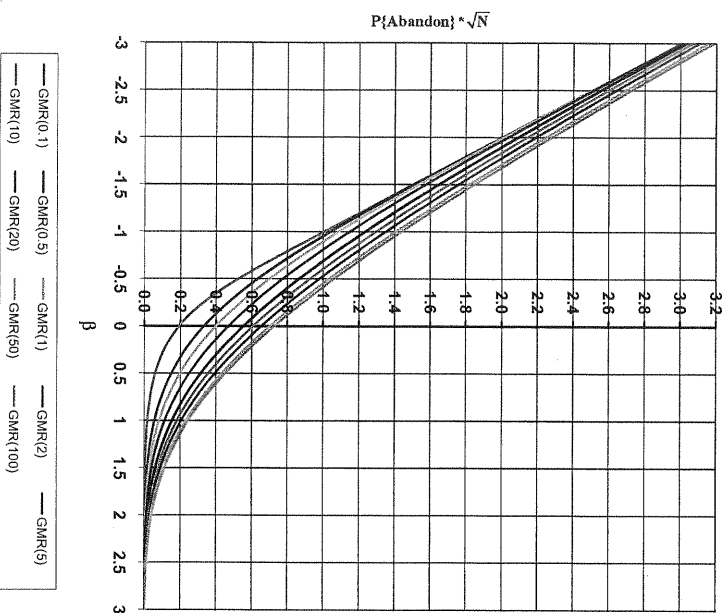
Understanding the Garnett Functions



- **Fix a staffing-level** (service-grade) and let patience \uparrow : then delays \uparrow ; in particular, the Garnett functions \uparrow to the Halfin-Whitt function (infinite-patience).
- **Fix a target delay-probability** (service level): then, as impatience \uparrow , less servers (smaller service-grade) are required to achieve the target (convincing managers to use Erlang-A).
- With $\beta = 0$ ($n = R$) and $\mu = \theta$, 50% are served immediately. Compare with Erlang-C in which $n = R + 0.5\sqrt{R}$ was required. But there is **no free lunch**: 2% abandon! (under $n = 400$) see next page.

Erlang-A: % Abandonment

$\%Ab \times \sqrt{n}$ vs. β , for varying (im)patience (θ/μ):



Note the behavior: slope $-\beta$, for (relatively) large negative β and over all (im)patience levels. For an explanation, think **ED**: $n = R + \beta\sqrt{R} = R - \gamma R$, hence $\gamma \approx -\beta/\sqrt{R} \approx -\beta/\sqrt{n}$, and γ is $P\{Ab\}$ in the ED-Regime.

7

“The Right Answer for the Wrong Reason” - Revisited

If $\beta = 0$, the QED staffing level $n \approx R + \beta\sqrt{R}$ becomes

$$n = R = \frac{\lambda}{\mu} = \lambda \cdot E[S],$$

which is equivalent to the following **deterministic** rule:

Assign a number of agents that equals the offered load.
(Common in stochastic-ignorant operations.)

Erlang-C: queue “explodes”.

Erlang-A: Assume $\mu = \theta$. Then $P\{W_q = 0\} \approx 50\%$.

If $n = 100$, $P\{Ab\} \approx 4\%$ (twice the value 2% in the graph - why?), and $E[W_q] \approx 0.04 \cdot E[S]$ (why?).

Overall, reasonable (good?) service level, which will in fact improve with scale. For example, with $n = 400$, both $P\{Ab\}$ and $E[W_q]$ reduce to half their value under $n = 100$ (why?).

(Note: Changes in n go hand in hand with same changes in λ , assuming μ remains fixed.)

The Effect of Patience:

Suppose now $\mu = 0.1 \cdot \theta$ (highly impatient customers).

Via the Garnett Functions, suffices $n = R - \sqrt{R}$ to achieve $P\{W_q = 0\} \approx 50\%$, but this comes at the cost of somewhat over 10% abandoning, with $n = 100$ (and 5% with $n = 400$); though $E[W_q]$ decreases to one fourth of the above, assuming μ remains unchanged.

8

Erlang-A in the QED Regime: Operational Performance Measures

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

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

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

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

$$P\left\{ \frac{W_q}{E[S]} > \frac{t}{\sqrt{n}} \mid W_q > 0 \right\} \approx \frac{\bar{\Phi}\left(\hat{\beta} + \sqrt{\frac{\theta}{\mu}} \cdot t\right)}{\bar{\Phi}(\hat{\beta})}$$

$$P\left\{ Ab \mid \frac{W_q}{E[S]} > \frac{t}{\sqrt{n}} \right\} \approx \frac{1}{\sqrt{n}} \cdot \sqrt{\frac{\theta}{\mu}} \cdot \left[h\left(\hat{\beta} + t \sqrt{\frac{\theta}{\mu}}\right) - \hat{\beta} \right]$$

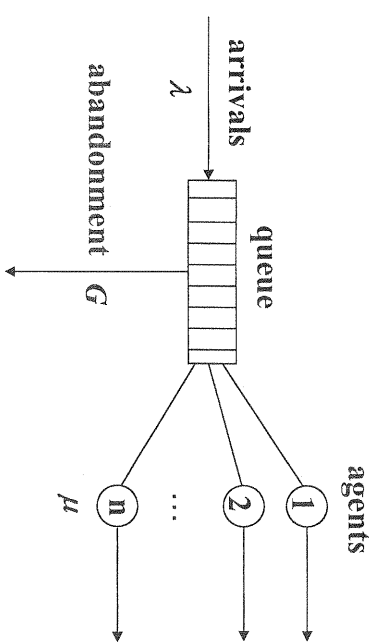
$$E\left[\frac{W_q}{E[S]} \mid Ab \right] \approx \frac{1}{\sqrt{n}} \cdot \frac{1}{2} \sqrt{\frac{\mu}{\theta}} \cdot \left[\frac{1}{h(\hat{\beta}) - \hat{\beta}} - \hat{\beta} \right]$$

Here

$$\bar{\Phi}(x) = 1 - \Phi(x),$$

$$h(x) = \phi(x)/\bar{\Phi}(x), \quad \text{hazard rate of } N(0, 1).$$

M/M/n+G in the QED Regime



Density of (im)patience G : $g = \{g(x), x \geq 0\}$.

Assume $g_0 \triangleq g(0) > 0$.

QED regime: $n \approx R + \beta\sqrt{R}$.

QED approximations: Use the Erlang-A formulae (from the previous page), substituting g_0 instead of θ .

How to estimate g_0 ? As $\hat{\theta}$ in Erlang-A!

Why? Recall **Erlang-A**: $P\{Ab\} = \theta \cdot E[W_q]$ used for estimating θ (either via $\hat{\theta} = [\#Abandoning] / [\text{Total Waiting Time}]$; or by regression of half-hours' $[\%Abandoning]$ over $[\text{Expected-Waits}]$).

M/M/n+G: It turns out that, in the QED regime:

$$P\{Ab\} \approx g_0 \cdot E[W_q].$$

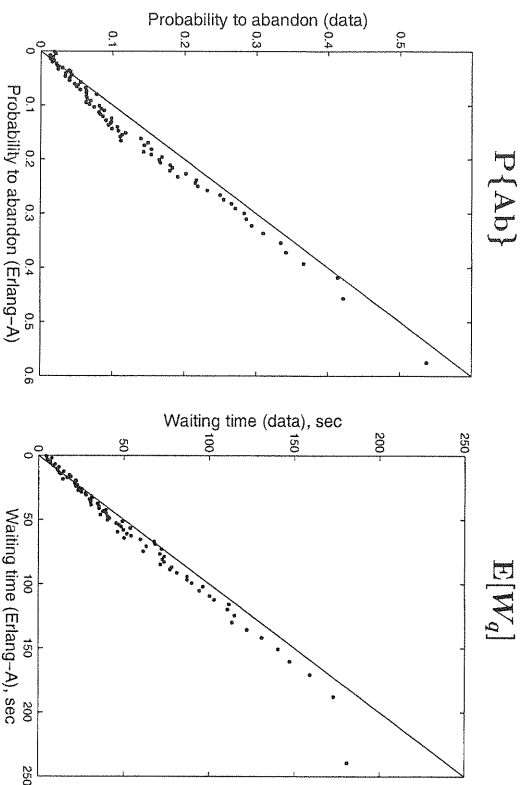
Hence, one estimates g_0 exactly as $\hat{\theta}$ in Erlang-A.

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

Question: Can one usefully apply the Erlang-A model to systems with non-exponential patience?

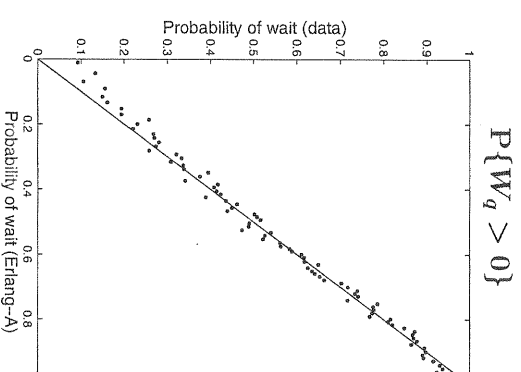
YES!

Erlang-A Formulae vs. Data Averages (Israeli Bank)



11

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



Summary:

- Points: Hourly data (averages) vs. Erlang-A predictions;
- Formulae with continuous n (special-functions) used to account for non-integer n ;
- Patience estimated via $P\{Ab\}/E[W_q]$;
- Erlang-A estimates provide close upper bounds.

12

Efficiency-Driven M/M/n+G (ED)

Let γ be a QOS parameter, $0 < \gamma < 1$.

Assume $G(x) = \gamma$ has a unique solution $x^* = G^{-1}(\gamma)$, at which $g(x^*) > 0$.

Staffing level:

$$n \approx R \cdot (1 - \gamma), \quad \gamma > 0.$$

$$\bullet \text{P}\{W_q > 0\} \approx 1.$$

• Abandonment-Probability converges to:

$$\text{P}\{\text{Ab}\} \approx \gamma \approx 1 - \frac{1}{\rho}.$$

• Offered-Wait converges to x^* :

$$\text{E}[V] \approx x^*, \quad V \xrightarrow{p} x^*.$$

• Waiting distribution (asymptotically):

$$W_q \xrightarrow{w} G^*, \quad \text{E}[W_q] \rightarrow \text{E}[\min(x^*, \tau)],$$

where G^* is the distribution of $\min(x^*, \tau)$, namely

$$G^*(x) = \begin{cases} G(x), & x \leq x^* \\ 1, & x > x^* \end{cases}.$$

Operational Regimes: Rules-of-Thumb

Assume that the **Offered-Load** R is not too small (more than several 10's for QED, more than 100 for ED and QD).

ED regime: $n \approx R - \delta R$, $0.1 \leq \delta \leq 0.25$.

- Essentially **all** customers are delayed;
- %Abandoned $\approx \delta$ (10-25%);
- Average-wait ≈ 30 seconds - 2 minutes.

QD regime: $n \approx R + \gamma R$, $0.1 \leq \gamma \leq 0.25$.
Essentially **no** delays.

QED regime: $n \approx R + \beta \sqrt{R}$, $-1 \leq \beta \leq 1$.

- %Delayed between 25% and 75%;
- %Abandoned is 1-5%;
- Average wait is one-order less than average service-time (eg. seconds vs. minutes).

Class 13

QED Qs - Part II: Erlang-A.

QED Q's: Economies of Scale; Staffing Moderate-to-Large Service Operations

Erlang-A (Abandonment); QEDing Palm's model;

Optional Reading Assignment: Read the article "Healthcare Call Centers: A Technology Migration", by Howard Burnett. Pay special attention to the following:

1. The Call Center Maturity Model (Figure 1);
2. The calls flow within the call centers (Figure 2), starting with the IVR, through the triage nurse, then to one of the advice groups (internal, pediatrics, obstetrics and gynecology) or appointment agents;
3. In regard to the last paragraph of the article, recall the article <http://iew3.technion.ac.il/serveng/Lectures/Retail.pdf> in which an attempt was made to define the Industrial Engineer of the Future.

Recitation 14: Introduction to operation regime. Shift scheduling.

HW 11: Empirical Analysis of a call Center via SEESat (and the Offered-Load).

This HW is based on real data, which you will be analyzing with SEESat. You will first identify problems with the operation of a call center, and then you will find staffing remedies for the difficulties found. The latter will require the use of 4CC.

A central role in the homework is the notion of Offered-Load, especially its time-varying version. There is also a part of the homework where you will check the validity of some congestion laws that were studied in class.

The due date for HW 11 is August 5, 2012.

Final Exam: $\leq 50\%$ of the Final Grade. Its structure, is as follows:

Question 1: From previous exams (see our WebSite), or from Recitations, or from Lectures, of very similar to one of these.

Question 2: From Homeworks.

[Those who know well the material from Lectures, Recitations, Homeworks, are very likely to get a final grade of at least 75-80.]

Question 3: A "Practical Question" with some theoretical insight. [With Question 3 answered well, one can get to a final grade of at least 85-90.]

Question 4: A "Theoretical Question" that requires deeper understanding. [Answering well Question 4 is required in order to get to the levels of 95-100 final grade.]

Topics that were left out, or just touched on:

- Skills-Based Routing
- Queueing Networks: Jackson and Non-Parametric
- New-Service Development (or Service-Engineering in Germany)
- Internet-based services (or Contact Centers)
- Appointments - managing demand (Hall, Chapter 8)
- Service Quality
- Forecasting/smoothing (F&F, Chapter 14)
- Location and (functional) design of service facilities (F&F, Chapters 6, 7)
- Marketing (Lovelock, who also has a book dedicated to the subject)
- Human resource management (Lovelock)
- Technology; Automation
- Convergence of Service and Manufacturing:
Field service, preventive maintenance, supply chain, life-time value
- Significant Service Sectors:
 - Health, Hospitality (Tourism), Financial, Transportation,
 - Telecommunication, Education, Professional Services,...

What's next?

- The "New-Age Industrial Engineer":
 - Industrial Engineers in Services: Banks, Hospitals, Government, etc.
 - Industrial Engineers in the the interface "Manufacturing – Services":
 - Consulting
 - Startups
- Research: Graduate Programs (Technion, Abroad); M.SC., Ph.D., TA'ing.
- Teaching:
 - Projects (Practical but Theoretically-Based)
 - Further Courses: Deeper (Q-Theory, Stochastic Processes); Broader (CRM, HRM, IE)

**PLEASE STAY IN TOUCH, ESPECIALLY IF YOU WORK IN "SERVICES"
OR ITS RELATIVES (MY GUESS - WITH PROBABILITY 0.75.)**

