

Class 7

~~Arrivals: Forecasting, and some loose ends~~ Service Times; Phase-type Distributions

Arrivals: Review

- Poisson processes: review;
- Forecasting arrivals;
- The Offered Load.

Defining, Modelling and Designing Service Times

- What is "Service-Time"? via Empirical analysis of face-to-face, telephone services; hospitals, ...
- Service time is a Statistical Distribution: lognormal, exponential.
- Service time is a Process: Phase-type distributions.
- Beyond Means and Beyond CV's.
- Stochastic Ordering.
- Subtleties.

Laws of Congestion: Old and New

The 0-th Law for (The) *Causes of Operational Queues* :

Scarce Resources and Synchronization Gaps (in DS-Project Networks);

The First Law of *Conservation* :

Little's Law for Customers, Service-providers and Managers.

Little's Law for the Offered Load (Utilization Profiles).

The Second Law of Completely *Random Arrivals* :

Levy/Watanabe Axioms of Randomness;

The Law of Poisson-Counting (Law of Rare Events);

The Law of Independent Memoryless (Exponential) Inter-arrivals;

The Brownian-Law of Rescaling & Centering Arrivals;

The Laws of Decomposition-Superposition.

The Third Law of *Human Service-durations* :

The Law of Phase-types for the Durations of Human Upaced Services;

The Empirical Law of Exponential/Log-Normal Durations.

The Fourth Law of *Sampling* :

Random Sampling: Wolff's PASTA = Poisson Arrivals See Time Averages;

Biased Sampling: Costs of Randomness; (Coefficient of Variation, or Form Factor).

Recitation 7. Statistical analysis of an arrival process.

Service Engineering

Class 7

Service Times (Durations, Processes)

Why Significant? eg. +1 second of 1000 agents costs \$500K yearly.

Why Interesting? Must accurately

Model, Estimate, Predict, Analyze, Design:

- Resolution: Sec's (phone)? min's (email)? hr's (hospital)
- Parameter, Distribution (Static) or Process (Dynamic)?
- Does it include after-call work?
- Does it include interruptions?
 - Whisper time, hold time, phones during face-to-face,...
- Does it account for return services?

How affected by covariates? How affects performance?

- Experience and Skill of agents (Learning Curve)
- Type of Customer: Service Type, VIP Status
- Time-of-Day: Congestion-Level
- Human Factor: Incentives, pending workload, fatigue
- Heavy-Traffic: The ED and QED Operational Regimes (later)

How to calculate Offered-Load? (towards Staffing)

Contents: **Service Times; Phase-Type Durations.**

- Service duration = **Statistical Distribution**:
 - **Empirical**: Histogram, CDF, Hazard Rates (Later);
 - **Parametric**: LogNormal, Exponential, Others.
- Empirical Introduction, mainly via DataMOCCA.
- Motivating Examples.
 - Designing an IVR/VRU.
 - Pooling a Service Network.
 - Long-term Care of the Elderly.
- Sample size.
- **What is Service Time (Duration)?**
A complex answer to a “simple” question:
 - Single vs. multiple visits.
 - After-Call Work (ACW); Utilization Profiles.
 - Time- vs. State-dependency.
 - Incentives (Call Center, Hospital)
 - Averages do not tell the whole story: the need for **Distributions**.
- Stochastic Ordering (of distributions).
- Service = **Stochastic Process**: Phase-type MJP.
- “Sufficient Statistics” in Heavy Traffic: ED, QED (later)
- Offered-Load (Work)

Parametric Distribution of Service Times

Most common parametric distributions in service systems are Exponential and Lognormal.

Exponential Distribution:

Density: $f(x) = \lambda e^{-\lambda x}$, $x \geq 0$,

Mean: $E[X] = \lambda^{-1}$,

Variance: $Var(X) = \lambda^{-2}$,

Coefficient of Variance: $C_v = \frac{SDV(X)}{E[X]} = 1$,

Median: $\lambda^{-1} \ln 2$.

An important property of the exponential distribution is that it is memoryless. This means that if a random variable T is exponentially distributed, its conditional probability obeys $\Pr(T > s + t \mid T > s) = \Pr(T > t)$ for all $s, t \geq 0$.

Lognormal Distribution:

Definition: X is a lognormal random variable if $\ln(X)$ is normally distributed with mean μ and variance σ^2 .

Density: $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\ln x - \mu)^2}{2\sigma^2}}$, $x \geq 0$,

Mean: $E[X] = e^{\mu + \sigma^2/2}$,

Variance: $Var(X) = e^{\mu + \sigma^2/2} (e^{\sigma^2} - 1)$,

Coefficient of Variance: $C_v = \sqrt{e^{\sigma^2} - 1}$,

Note that CV does not depend on μ . For small σ ($\sigma < 0.5$), one can use $CV \approx \sigma$.

Median: e^μ .

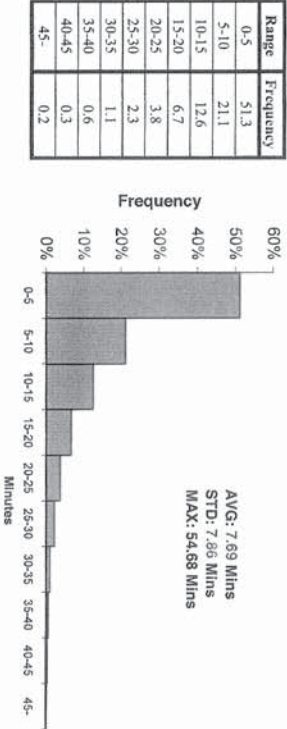
Local Municipalities Service Time

Department	Station No.	Total Customers	Avg. Arrival Rate (1/Hr)	Avg. Service Time (Mins)	STD (Mins)	Maximal Service Time (Mins)	Utilization	Avg. Waiting Time (Mins)
Water	N/A	187	1.8 ± 0.2	8.87 ± 1.0	8.15	54.68	13.3%	4.76
Tellers	N/A	1328	12.6 ± 0.5	8.82 ± 0.4	8.55	49.37	30.8%	7.73
Cashier	N/A	757	7.2 ± 0.4	6.64 ± 0.4	6.94	29.95	79.7%	3.89
Manager	N/A	190	1.8 ± 0.2	7.99 ± 1.0	8.44	38.97	24.1%	9.16
Discounts	N/A	317	3.0 ± 0.3	4.59 ± 0.4	4.54	36.72	23.1%	3.65

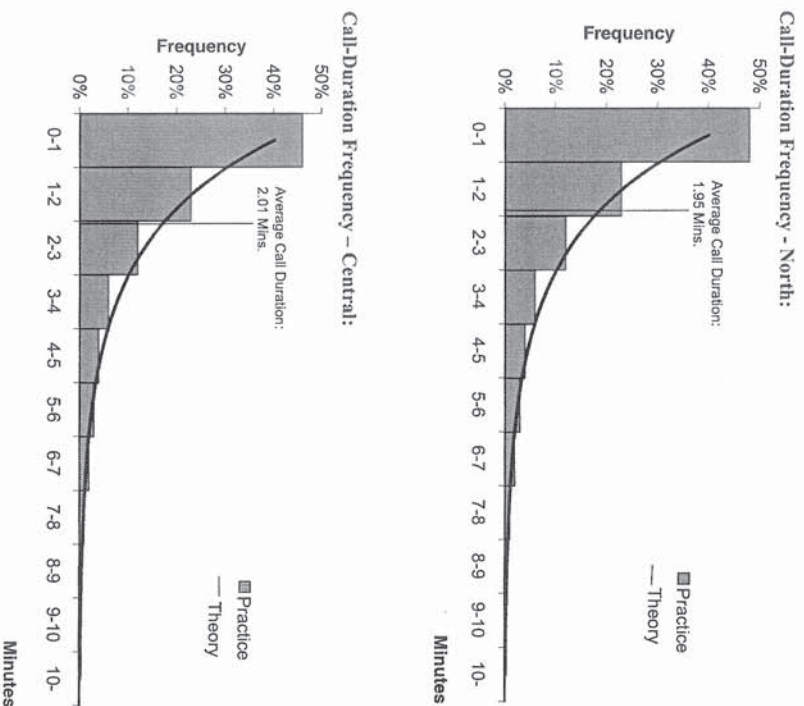
Water	1	57	N/A	7.80 ± 1.70	7.61	31.28	6.5%	N/A
	2	130	N/A	9.34 ± 1.20	8.37	54.68	19.3%	N/A
	3	336	N/A	9.04 ± 0.80	8.93	49.05	48.2%	N/A
	4	208	N/A	9.93 ± 1.00	8.82	49.12	33.0%	N/A
	5	417	N/A	8.97 ± 0.70	8.55	49.37	59.4%	N/A
Tellers	6	144	N/A	9.53 ± 1.20	8.75	41.70	21.8%	N/A
	7	156	N/A	8.03 ± 1.10	7.96	35.27	19.8%	N/A
	8	67	N/A	3.74 ± 0.70	3.58	21.03	4.0%	N/A
Cashier	9	757	N/A	6.64 ± 0.40	6.94	29.95	79.7%	N/A
Manager	10	190	N/A	1.99 ± 1.00	8.44	38.97	24.1%	N/A
Discounts	11	317	N/A	4.59 ± 0.40	4.54	36.72	23.1%	N/A

*Service time ranges given with 90% confidence.

Service Time Histogram – Overall:



Service Times: Exponential (Phone Calls)



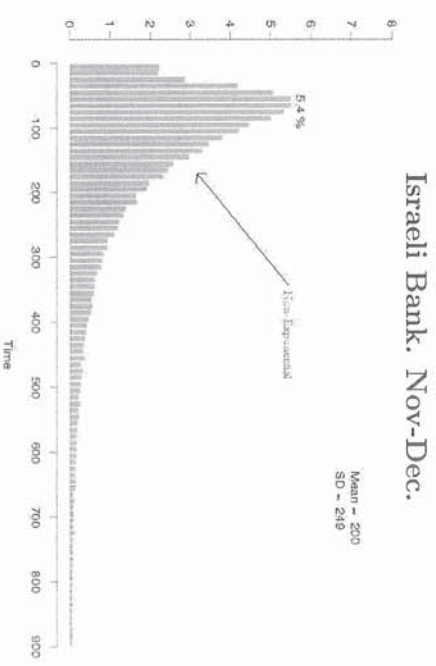
Q: How to recognize “Exponential” when you “see” one?

A: Geometric Approximation

5

LogNormal Distribution

Empirically prevalent in call centers (overall, service types, individual agents), but yet **no** theoretical explanation.



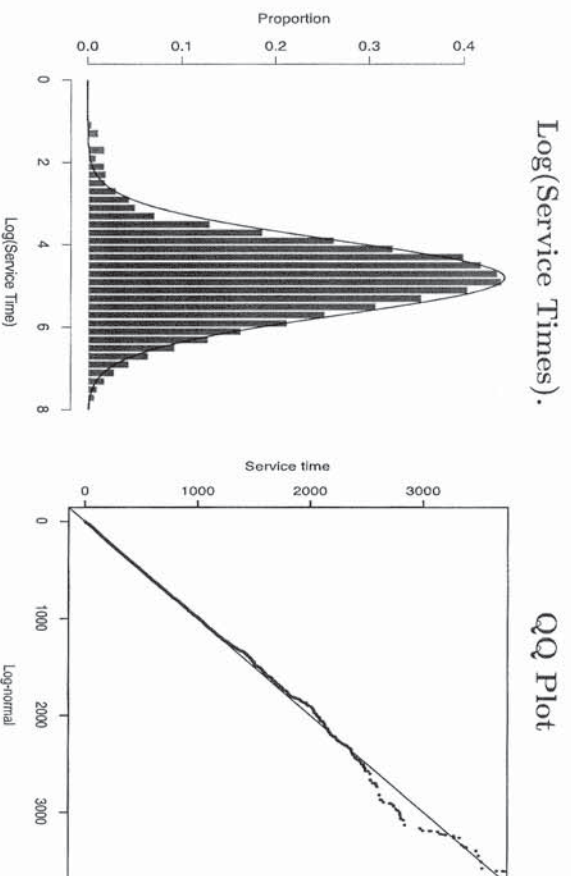
Good in **statistical** models
(eg. regression of $\log(\text{service-time})$).

Not so good for **queueing** models
(which typically “prefer” Exponential durations).

The practical good news for service time distribution in **queueing** models
CV is more important, if tail of distribution is similar one can use exponential distribution.

6

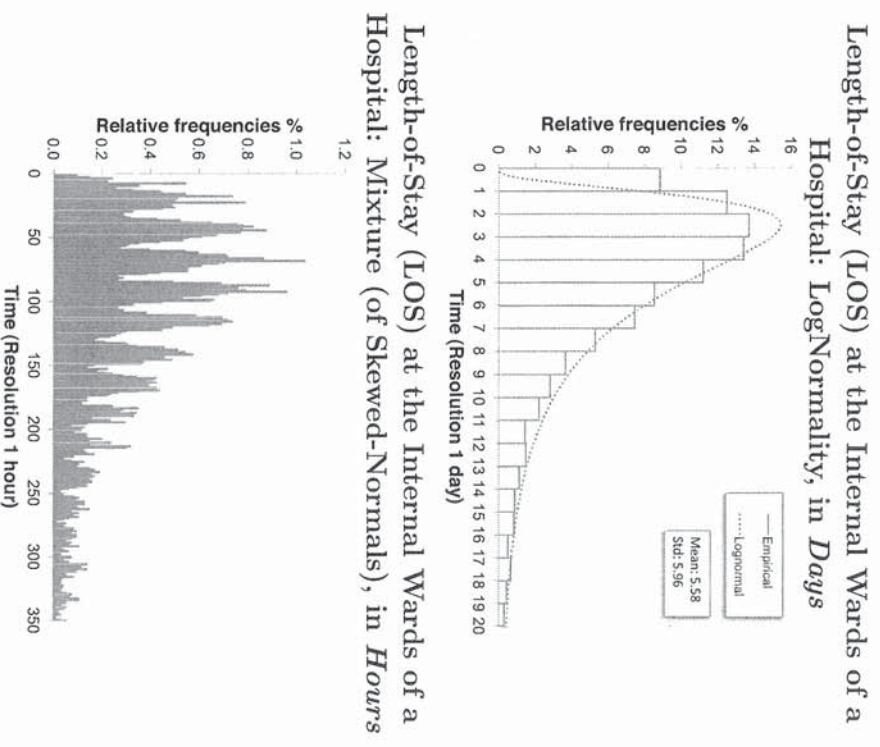
Validating LogNormality of Service Times



QQ Plots will be reviewed at the Recitation

7

Beyond Standard Distributions: Service Times at Hospital



8

Service Times: Mixture Model

A mixture model represents the presence of sub-populations within an overall population.

Finite mixture: Given a finite set of probability density functions $f_1(x), \dots, f_n(x)$, and weights w_1, \dots, w_n such that $w_i \geq 0$ and $\sum_{i=1}^n w_i = 1$, the mixture distribution can be represented by writing the density, f , as a sum (which is a convex combination):

$$f(x) = \sum_{i=1}^n w_i f_i(x).$$

Moments:

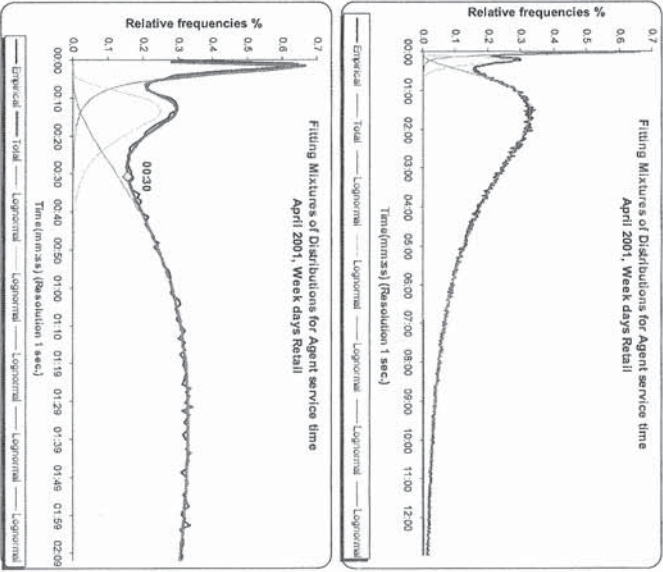
$$E[X] = \mu = \sum_{i=1}^n w_i \mu_i,$$

$$E[(X - \mu)^2] = \sigma^2 = \sum_{i=1}^n w_i (\mu_i^2 + \sigma_i^2) - \mu^2.$$

Example: Service time distribution in a call center, Length-of-stay in Maternity Wards.

Service Times: Mixture

LOS in Call Center as a mixture of LogNormals

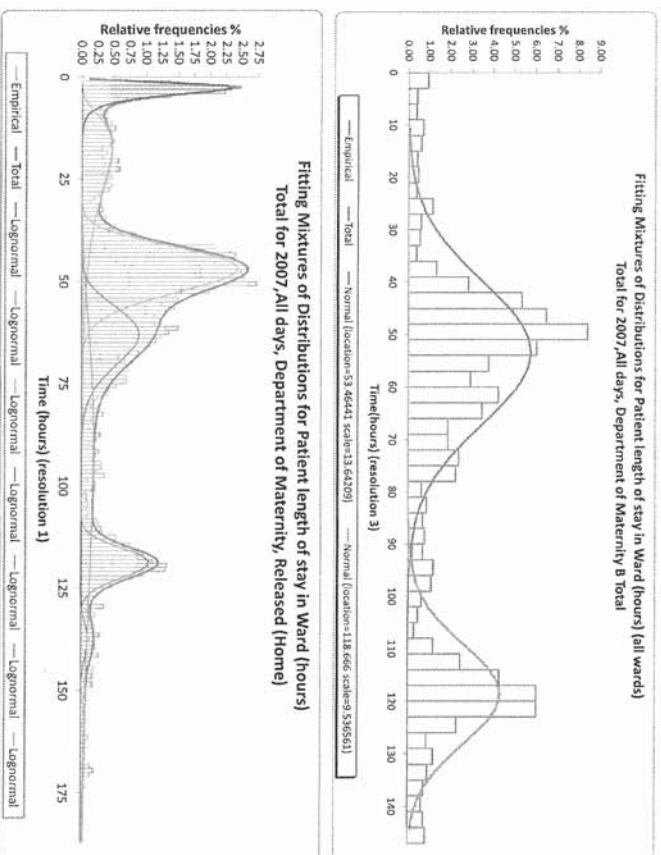


Going to the previous Excel Sheet, to view the corresponding Tables (by scrolling it down), one notes that the main component has weight 91% in the mixture – its role in the chart is to fit the part beyond 30 seconds, which it does very well.

Parameters Estimation	
Components	Mixing Proportions
1. Lognormal	3.08
2. Lognormal	3.58
3. Lognormal	91.14
4. Lognormal	1.86
5. Lognormal	0.34

Service Times: Mixture

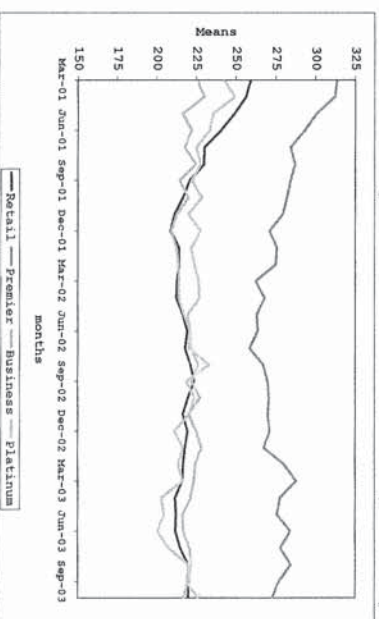
A Mixture Distribution for LOS in Maternity Wards



11

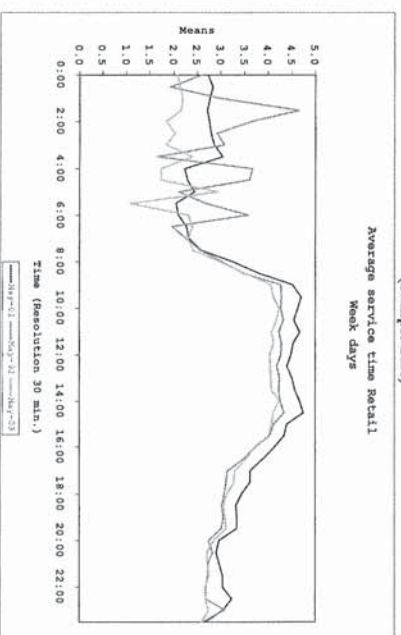
Service Times: Trends and Stability

USBank Average Customer Service Time, Weekdays



USBank Average Customer Service Time, Telesales

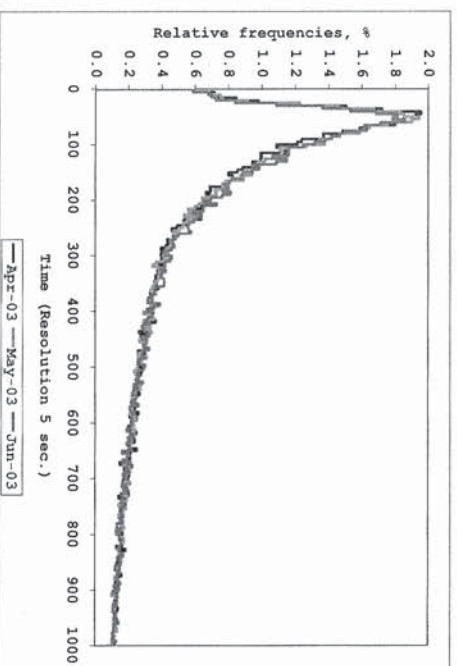
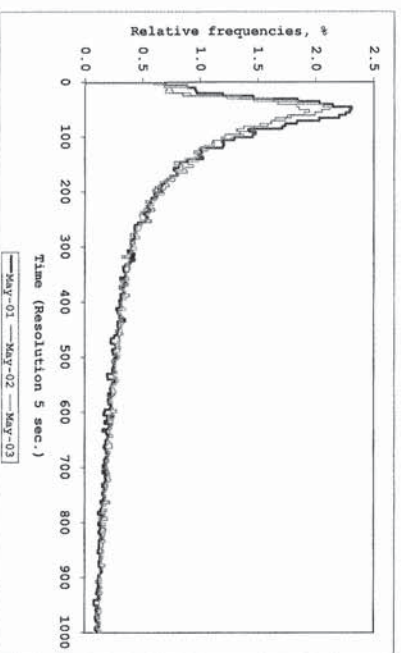
US Bank: Dynamics of average customer service time for Retail calls
(Sample Size)



12

Service Times: 5 Sec's Resolution

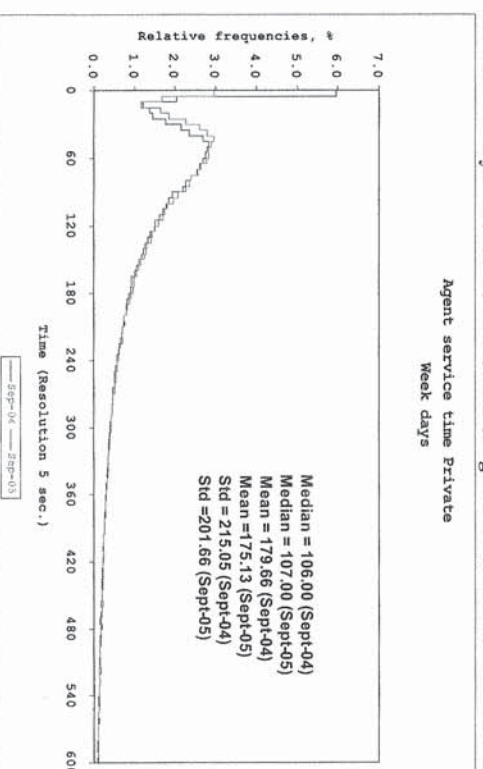
USBank, Service-Time Histograms for Telesales (MOCCA)



13

Service Times in Israeli Telecom

IL Telecom: Dynamics of the distribution of agent service time for Private calls



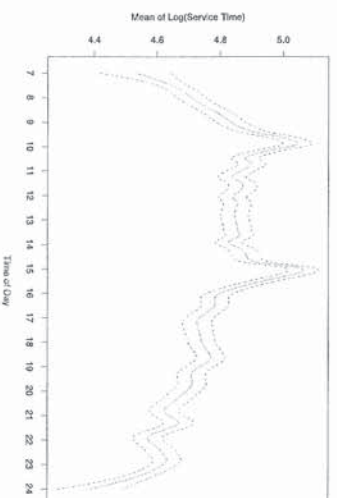
- Overall pattern seems close to LogNormal (except for the very short service times);
- Histograms of different months are **very** similar;
- Reason for short service durations unknown here.

14

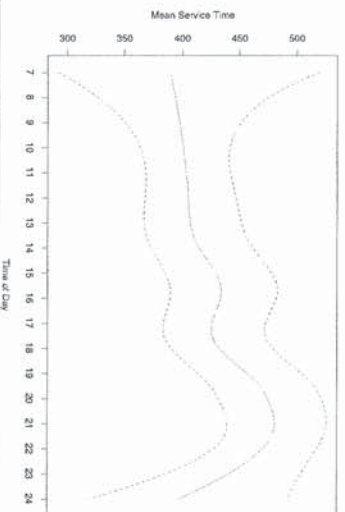
Service Times: The Human Factor, or Why Longest During Peak Loads?

Mean Service Time vs. Time-of-Day

Regular Service (PS)

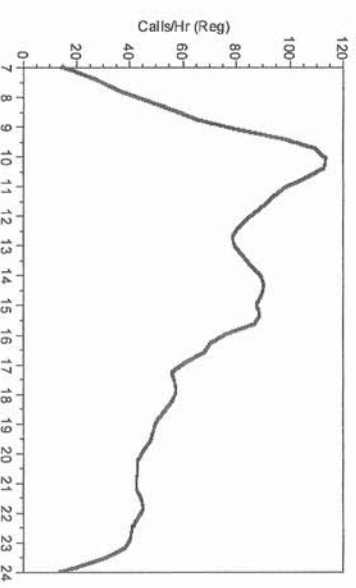


Internet Consulting (IN)



Service Times vs. Arrival Rates

Regular Service (PS): Arrival Rate



At 10:00 & 15:00: longest services and peak arrival rates?

Possible Reasons:

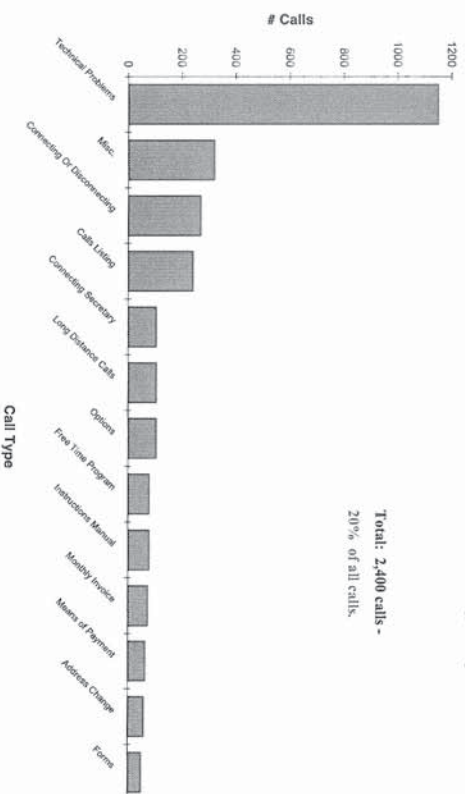
1. *Services are longer* during congestion since customers start with complaints.)
2. *Agents are slower* at times of peak loads.
3. Customers that arrive during peak hours require, for some reason, *longer service*.
4. An additional (*human*) reason will be provided after we study *customers' impatience*.

Service Time \neq Contact-Time

Common (Often Too Common):

- Customers routed for additional services (vs. “First-Time-Resolution”);
- Servers interrupt face-to-face service with a phone-call (vs. the increasingly prevalent “Central Call Center”);
- Agents place customers on hold, eg. technical consultation with veterans;
- Agents can be engaged in non-phone activities, eg. ACW Time (After-Call Work).

Reasons for Redials in a Cable Company



17

Calculating (Mean) Service Time

First approach:

Sum up components of the “service time”, then add related activities of servers.

Second approach (Avoids Ambiguities):

Fix a time interval (eg. a shift).

$$\text{Mean Service Time} = \frac{\text{Available Time} - \text{Idle Time}}{\text{Number of Calls}},$$

where

Available Time = # Agents \times Interval Duration, and

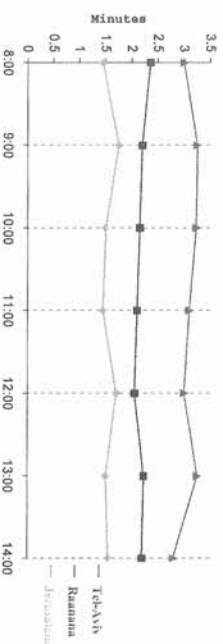
Idle Time is summed over all agents.

20

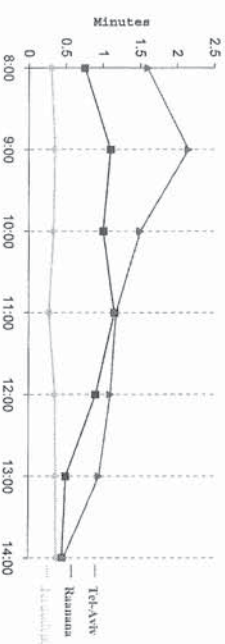
Israel Electric Company: 3 Centers

Service Performance

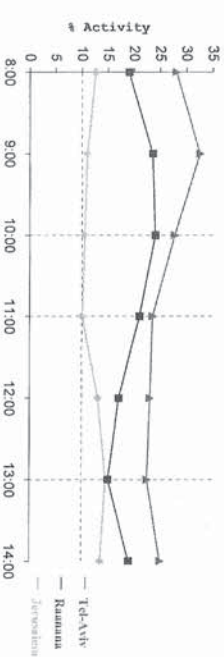
Service Time - Average:



Waiting Time - Average:



% Abandonment:



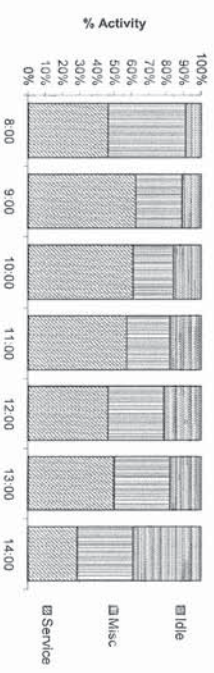
18

Israel Electric Company: 3 Centers

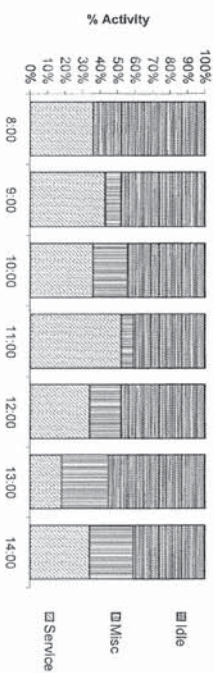
What is "Service Time"?

Utilization Profile in 3 Call Centers Doing the Same Thing

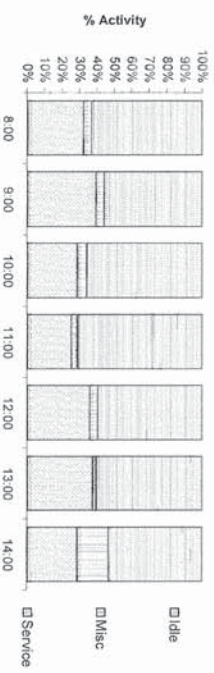
Tel-Aviv:



Ramatana:

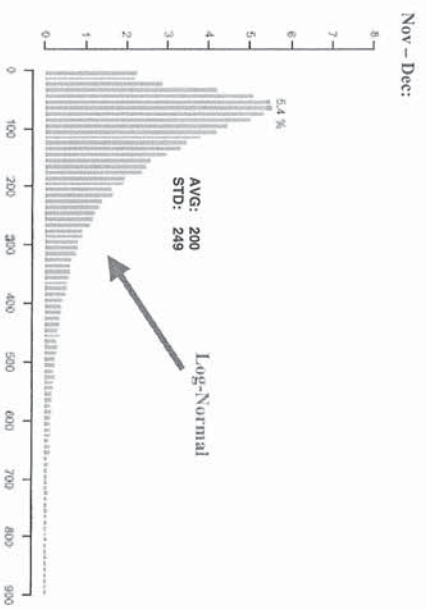
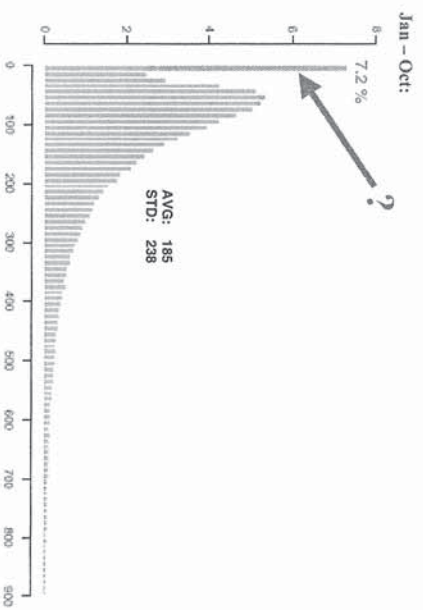


Jerusalem:



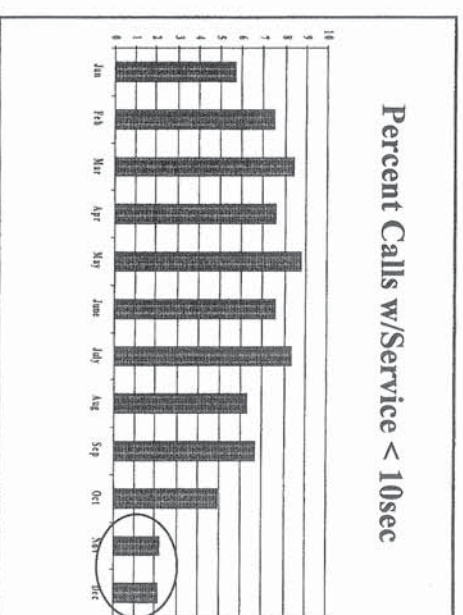
19

Beyond Data Averages Short Service Times



21

Short Service Times: Time Series



22

Short Service Times: Individuals

Mandelbaum, Sakov and Zeltyn

52

Table 52: Number of calls handled by an agent

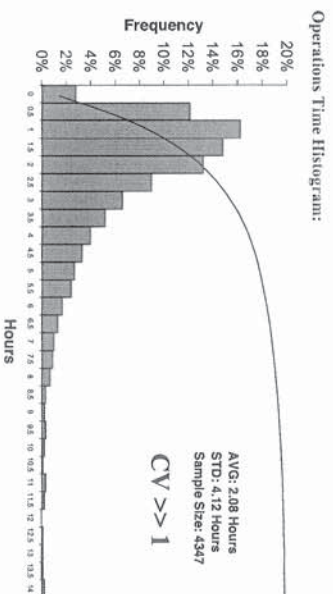
	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
AVT	0	0	0	1117	2208	2019	2789	2710	1417	2028	2533	2535
AVNI	1493	1736	642	539	1786	2219	2092	2392	1136	1885	1888	2136
BASCH	999	1164	1708	1155	982	906	855	2185	1973	1055	1336	1242
BENSON	1283	1135	0	1053	1108	1016	1682	1298	1076	1303	1546	1176
DARLON	309	515	653	519	577	436	309	370	297	194	425	128
DOORT	696	1047	0	811	546	862	750	2228	1319	1384	1640	1605
ELI	387	508	777	447	560	436	395	458	416	363	502	362
GELBER	333	143	510	427	859	281	386	332	67	179	165	269
GILL	665	614	1155	803	1108	974	418	0	355	466	412	298
KAZAV	1995	1603	1240	1451	1731	2251	1737	1168	729	1570	1047	2038
MEIR	0	0	0	0	0	0	127	344	318	280	406	454
MOULAI	1360	1223	1591	1351	1866	1980	2416	2152	1526	1940	1793	515
PINHAS	79	40	359	244	31	311	422	241	143	105	51	63
ROTH	0	0	397	1292	1928	1967	1831	1749	1625	1914	1458	1038
SHARON	1985	1674	2780	1538	2563	2657	2537	2873	1803	1535	2532	2140
STEREN	0	1043	2294	1516	2163	2231	1423	2465	1672	709	2375	2565
TOVA	1923	1679	1562	1059	1464	1380	1890	1811	1361	1971	941	0
YICKY	895	0	0	0	1006	1378	1415	1674	1472	1582	1641	1990
YIFAZ	1312	1901	1745	1305	1464	1076	780	90	1137	1315	0	0
YITZ	1771	1791	1402	1303	1355	1367	1009	69	705	1743	2420	2353
ZOHARI	891	1144	1398	1148	1479	1450	980	1494	1423	1359	1504	1094
ZZARIE	0	0	0	0	0	0	0	56	225	315	432	534
ZZELNOR	0	0	0	0	0	0	0	45	352	288	222	310
ZZELYAL	0	0	0	0	0	0	0	95	331	428	579	618
ZZEFAT	0	0	0	0	0	0	0	94	260	314	215	0
ZZINOR	0	0	0	0	0	0	0	84	250	136	136	138
ZZINRUT	0	0	0	0	0	0	0	116	327	474	387	545
ZZOPREIZ	0	0	0	0	0	0	0	71	311	260	242	334
ZZSPIECEL	0	0	0	0	0	0	0	71	311	260	153	322

Table 53: Number of calls with short service time

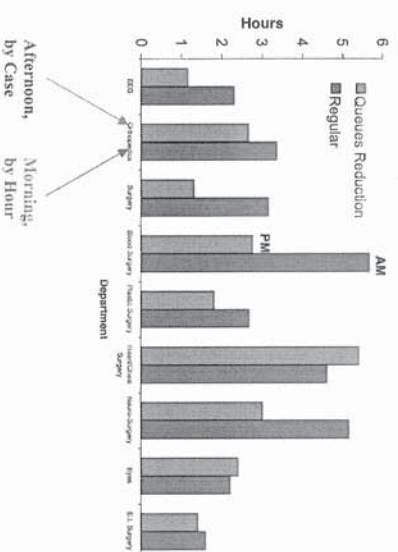
	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
MOULAI	235	230	356	369	614	695	597	490	455	4	1	
AVT	0	0	0	47	111	144	295	221	121	76	35	26
AVNI	11	13	4	4	5	6	25	16	18	4	8	11
DARLON	2	11	8	9	10	7	1	1	1	0	0	0
ELI	9	7	10	12	22	18	15	4	8	3	6	5
KAZAV	57	40	48	44	48	63	40	27	13	18	6	6
MEIR	0	0	0	0	0	0	1	8	3	1	2	0
PINHAS	3	0	58	25	4	14	11	6	8	1	0	0
ROTH	0	0	10	10	36	21	43	25	52	31	3	6
SHARON	58	49	86	52	67	78	66	63	38	23	43	49
TOVA	52	163	269	132	231	193	100	109	207	190	6	0
ZOHARI	4	8	12	22	17	20	9	14	5	7	10	7

23

Service Times: The Human Factor, or Even “Doctors” Can Manage



Operations Time - Morning vs. Afternoon:



Ethical?
Even Doctors Can Manage!

24

Human factors: Learning and Forgetting

Service Time Trends of Individual Servers

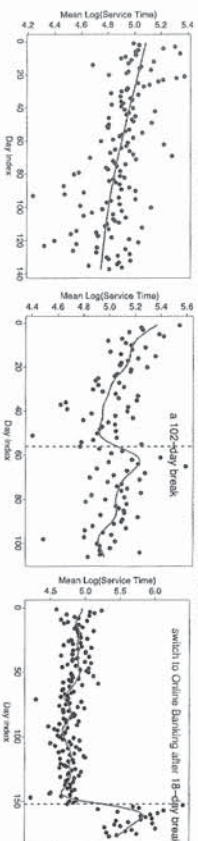


Figure 5: Long-term trend of daily average service time

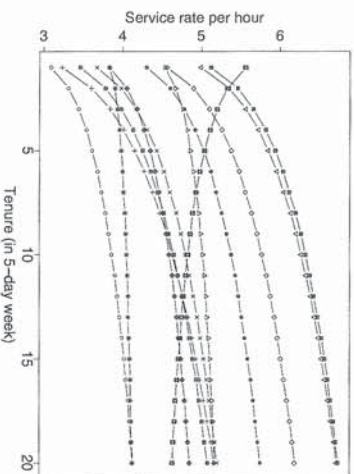


Figure 6: Daily learning curves of the 12 agents at site S

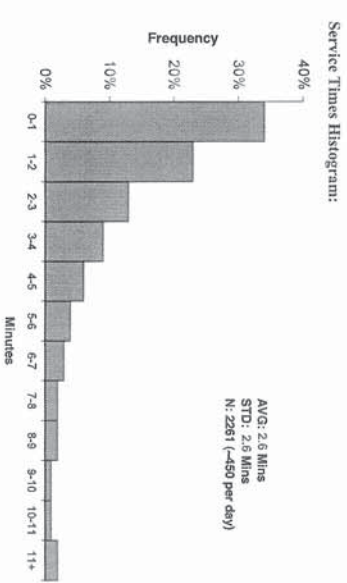
Classic learning models:

Assuming service times are lognormal distributed. y_{jk} - the service time of the k th call during the j th day. n_j - the total number of calls served by this agent during the j th day. Define $z_{jk} = \log(y_{jk})$. Then, the basic learning model is:

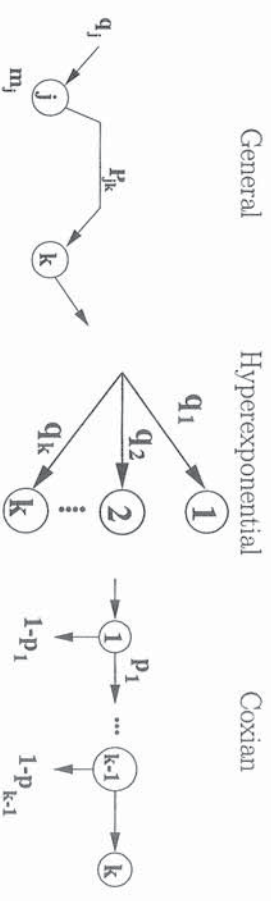
$$z_{jk} = a + b \log(j) + \epsilon_{jk}, \quad \epsilon_{jk} \sim N(0, \sigma_j^2)$$

Service Times: from Exponential to Phase-Type

Static Model: Exponential Duration Face-to-Face Services in a Government Office

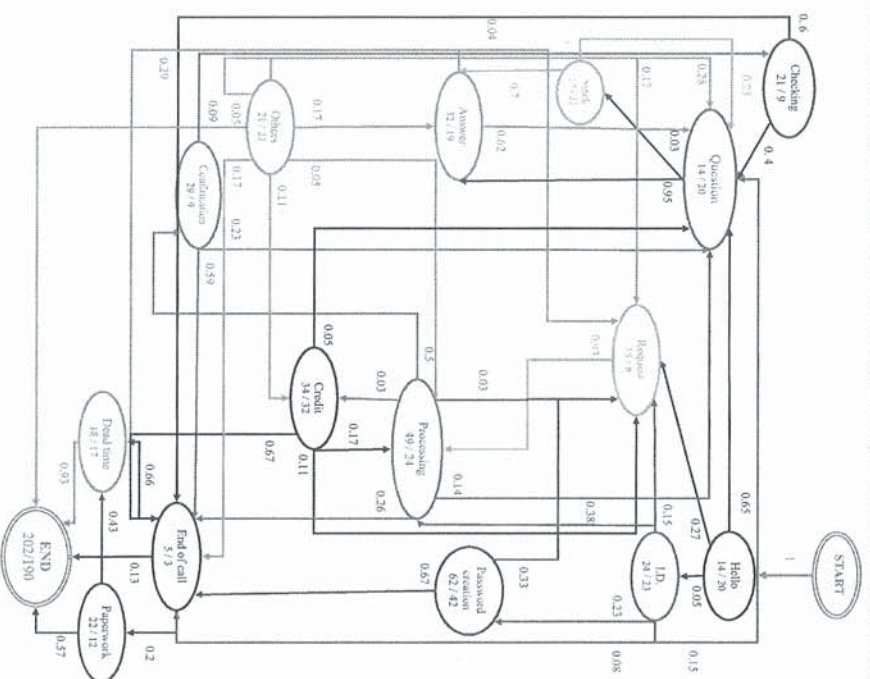


Dynamic Model: Phase-Type Duration



Phase-Type Model of a Telephone Call

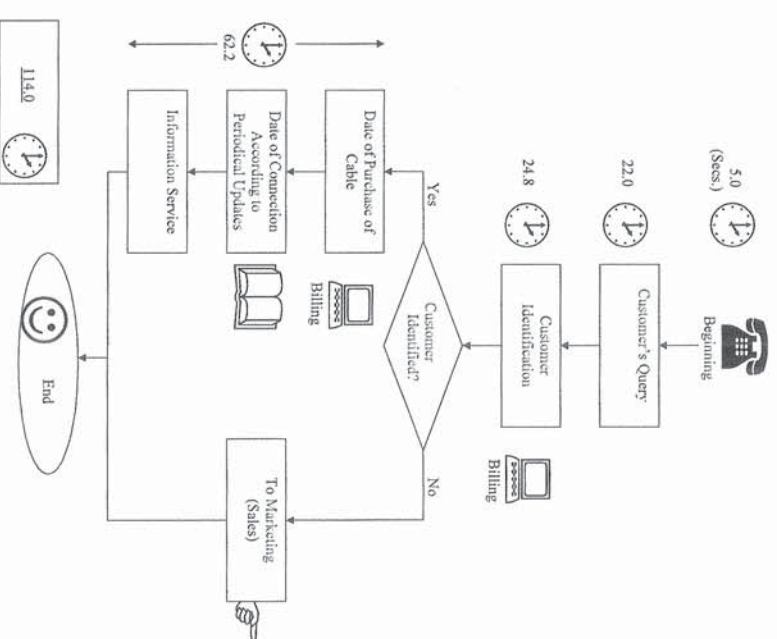
Figure 20: Phase-Type Model of a Telephone Call (# within phases: Mean/STD)



26

Service Times: Phase-Type Model

Late Connections



? Where does human-service start / end (recall 144)?
"Average" picture.

27

Phase-Type Service Times (Durations).

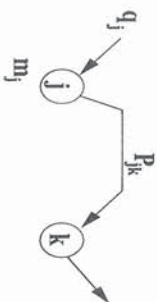
Service-Time = a sequence/collection of tasks, of an *exponential* duration.
There are K types of tasks, indexed by $k = 1, \dots, K$.

m_k = expected duration of task k ;

q_k = % of services in which k is first;

P_{jk} = % of incidences in which task j is immediately followed by k . $\bar{P} = [P_{jk}]$

$1 - \sum_{k=1}^K P_{kk}$ = probability to end service at k .



Fact: service = finite number of tasks $\Leftrightarrow \exists [I - P]^{-1}$
Indeed, $[I - P]_{jk}^{-1}$ = expected number of "visits to k ", given j was first.
 $(q[I - P]^{-1})_k$ = expected number of "visits to k ".

As will be articulated below, service-time duration is *Phase-type* (PH).
(Assuming independence among task-durations.)

Definition. Phase-type distribution = absorption time of a finite-space continuous-time Markov chain, with a single absorbing state.

Formally: $X = \{X_t, t \geq 0\}$ Markov on states $\{1, 2, \dots, K, \Delta\}$, with infinitesimal generator

$$Q = \begin{bmatrix} 1 & & & \\ \vdots & R & r & \\ K & 0 & \dots & 0 & \Delta \end{bmatrix} \quad \begin{array}{l} \bullet \Delta \text{ absorbing} \\ \bullet r = -R1 \\ \bullet 1, \dots, K \text{ transient} \end{array} \quad \begin{array}{l} (\text{since } q_{\Delta\Delta} = 0) \\ (\text{since } Q1 = 0) \\ \Leftrightarrow \exists R^{-1} \text{ (fact)} \end{array}$$

and initial distribution (of X_0) is given by $(q_1, \dots, q_K, 0) = (q, 0)$.

Recall:

$$P\{X_t = k\} = \sum_j q_j [\exp(tR)]_{jk} = q[\exp(tR)]_k$$

Define: $T = \inf\{t > 0 : X_t = \Delta\}$ has phase-type distribution, say $F_T(\cdot)$.

Claim: $F_T(t) = 1 - qe^{tR}1$, $t \geq 0$.

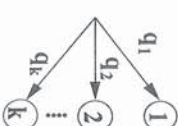
Proof. $P(T > t) = P\{X_t \neq \Delta\} = \sum_k q(e^{tR})_k = qe^{tR}1$.

Parameters:

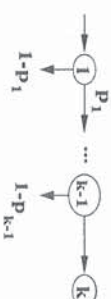
density $f_T(t) \equiv qe^{-tR}$
Laplace transform $\int_0^\infty e^{-st} f_T(dt) = q[sI - R]^{-1}r$
 n th moment $\int_0^\infty t^n f_T(dt) = (-1)^n n! qR^{-n}1$
(mean $= -qR^{-1}1$)

Special Cases:

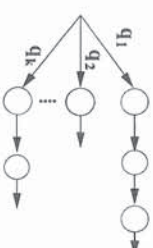
- Exponential (μ): $R = [-\mu]$ and $q = 1$.
- Erlang: $\rightarrow [1] \rightarrow [2] \rightarrow \dots \rightarrow [K]$ iid tasks / phases ($C^2(T) = \frac{1}{K}$).
- Generalized Erlang: exponential phases in series (tandem) ($C^2 < 1$).
- Hypereponential: K tasks in parallel (mixture) ($C^2 > 1$).



- Coxian: K phases; end at phase k with probability p_k .



- Minimum of exponential random variables is exponential.
- Max of exponential random variables is phase-type: e.g., $X_i \sim \exp(1)$ iid. This easily implies that $E(\max X_i) = \sum_{i=1}^n \frac{1}{i}$, $\text{Var}(\max X_i) = \sum_{i=1}^n \frac{1}{i^2}$ bounded!
- Erlang mixtures:



Importance of Phase-type distributions.

- Empirical + wishful thinking: homogeneous human tasks are exponential.
- Richness: the family of phase-type distributions is dense among all distributions on $[0, \infty)$. For every non-negative distribution G , there exists a sequence of phase-type distributions $F_n \ni F_n \Rightarrow G$. (In particular, we can guarantee convergence of any finite number of moments.)

Dense subfamilies: Coxian, Erlang mixtures.

For Erlang mixtures, this can be explained by the following two facts:

1. The family of discrete distributions is dense.
2. Constants can be approximated by Erlang distributions. Therefore, discrete distributions can be approximated by Erlang mixtures.

- Modelling, via the *method of phases*. For example, consider M/PH/1 queue (see HW).

M/PH/1: state-space is (i, k) (i = number in queue; k = phase of service) or 0 ; $e, \delta, 0 \xrightarrow{M_k} (1, k)$.

Representation directly in terms of (q, R, m) .

Denote here $R = [I - P]^{-1}$ (as in Mandelbaum & Reiman).

Average work content $E(T) = qRm$ ($= \sum_j q_j R_{jj} m_k$).

Moments: $E(T^n) = n! q(RM)^n q$, where $M = \begin{bmatrix} m_1 & & 0 \\ & \ddots & \\ 0 & & m_K \end{bmatrix}$

$$\frac{E(T^2)}{2(E(T))^2} = \frac{1 + C^2(T)}{2} = \frac{q(RM)^2 1}{(qRM)^2}$$

Pooling Services: Municipality Service System

Current state:

Service Station	No. of Work Stations	Arrivals per Hour	Occupancy	Average Waiting Time [Minutes]	Average Queue Length
Collection - Front Office	4	17.50	0.48	28.08	8.19
Collection - Immigrants	2	8.43	0.76	30.92	4.34
Collection - Back Office	6	3.20	0.13	13.45	0.72
Cashier	2	22.80	0.70	9.73	3.70
Assessment - Front Office	4	11.73	0.50	15.58	3.05
Assessment - Back Office	2	22.80	0.05	11.72	0.07
Land registry office	2-4	2.00			

Recommended staffing in overloaded periods (using model):

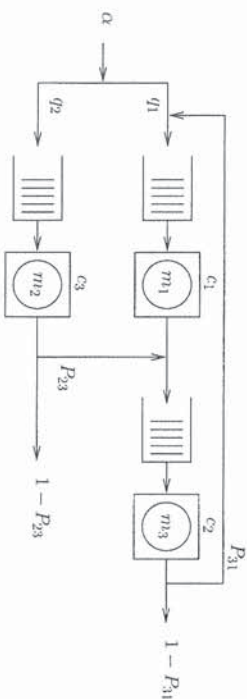
Service Station	Arrivals per Hour	Service Time	Recommend No. of Clerks	Percent of Customers Served	Average Waiting Time [Minutes]	Probability of Waiting	Spaces for Waiting
Collection - Front Office	23.40	6.98	6	45	2.1-4.8	0-0.06	4-6
Collection - Immigrants	4.50	14	3	35	7.1-16.2	0.006-0.1	3-5
Collection - Back Office	11.80	12	4	59	7.3-16.8	0.007-0.27	6-9
Cashier	31.40	3.5	3	61	3-6.6	0-0.37	7-10
Assessment - Front Office	16.00	10.9	6	48	3.53-8.1	0-0.09	5-7
Assessment - Back Office	0.60	18.18	2	9	10-22.8	0.015-0.027	2

Total of 24 clerks.

Pooling Services: Municipality Queueing Network (Server's Perspective)

Pooling Services: Municipality Service Times per Service Position

Figure 3 A Specialized Model with Task Repetition and Feedback



Station 1 - Collection;
 Station 2 - Assessment;
 Station 3 - Cashier.

Dept.	Server ID	Service Time Avg. (Min)	Std. Deviation	Utilization %	Service Time Max. (Min)	Total Services
Collection - Front Office	1	7.55 ± 0.68	7.96	37	79.32	370
	2	5.42 ± 0.33	6.27	68	105.20	951
	3	6.51 ± 0.50	6.94	44	63.33	510
	4	8.41 ± 0.75 ± 8.90	8.90	42	58.15	377
Collection - Immigrants	5	11.59 ± 0.80 ± 10.88	10.88	76	74.60	493
	6	10.32 ± 0.52	8.98	78	50.87	569
Collection - Back Office	7	10.80 ± 1.98	12.82	16	93.73	114
	8	9.07 ± 3.56	11.50	3	52.07	28
	9	18.32 ± 4.90	20.34	10	113.57	47
	10	23.39 ± 5.52	17.75	9	63.77	28
Cashier	11	11.99 ± 3.16	14.75	9	70.30	59
	12	16.73 ± 2.34	16.08	28	88.68	128
	13	2.51 ± 0.21	4.92	48	52.18	1460
	14	3.86 ± 0.18	4.16	72	46.92	1416
Assessment - Front Office	15	13.74 ± 1.07	12.02	62	69.68	340
	16	10.88 ± 0.92	10.60	52	87.92	363
	17	6.66 ± 0.50	6.68	42	49.93	473
Assessment - Back Office	18	11.22 ± 1.30	13.81	45	100.60	302
	19	19.29 ± 5.64	19.99	8	78.27	34
Total	20	12.2 ± 3.86	8.47	3	29.28	13
		7.24 ± 0.10	9.10			8075

- 90% confidence intervals
- 7364 distinct customers

Recall: Exponential $\Rightarrow E = \sigma$ (i.e. CV=1)

Pooling Services: Services Classification

A Classification of Service Tasks

Complexity

complex

short

long

• r1

• r2

• fs

• f1

Duration

Customer

Frequency Manager

rare

frequent

Improvement efforts

FC L
↓
RSS

Pooling Services: Municipality Server Recommendation

Recommendations: Flexible clerks for all activities. Change from Figure 3 to Figure 4.

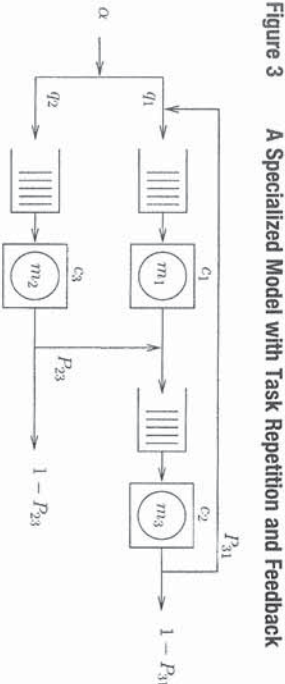
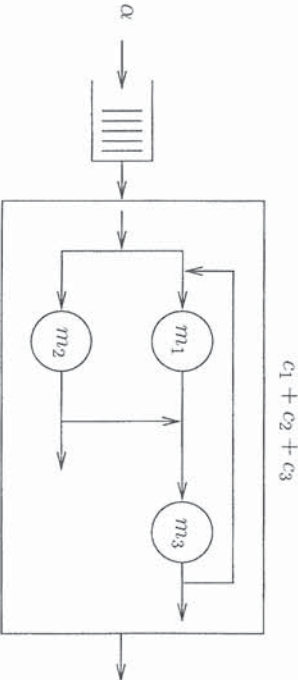


Figure 4 The Flexible Model, under Complete Pooling, that Corresponds to Figure 3



Pooling Services: Municipality Server Recommendation

State of system under recommendations:

Hour	Arrival Rate	Staffing	Occupancy	Waiting Room Size	Average Waiting Time [Minutes]	% Waiting More than 10 Min
7:30-8:30	36.3	(7) 8	(69) 60	8-12	3.20	4.7
8:30-9:30	79.4	(13) 14	(82) 76	14-22	3.10	3.6
9:30-10:30	87.4	(15) 15	(78) 78	16-24	3.05	3.9
10:30-11:30	85.4	(14) 15	(81) 76	15-22	2.85	3
11:30-12:30	64.5	(11) 12	(78) 72	12-18	3.00	3.7
12:30-1:30	24.5	(6) 7	(54) 46	6-8	2.70	2.7
1:30-2:30	24.2	(6) 7	(54) 46	5-8	2.70	2.6
2:30-3:30	30.6	(7) 8	(58) 51	6-9	2.71	2.2
3:30-4:30	11.3	(4) 5	(34) 30	4-5	2.65	3.4

- Number of work-stations: 15.
- Staffing change over time between 5 and 15.
- Guidance on matching available agent to waiting customer.
- Standardization of services and work procedures.
- Turnover clerks to achieve high occupancy.
- Possible separation between Russian speaking clerks.

On Pooling in Queuing Networks

Avishai Mandelbaum • Martin I. Reiman
Faculty of Industrial Engineering and Management, Technion, Haifa, Israel
Bell Labs, Lucent Technologies, Murray Hill, New Jersey 07974

We view each station in a Jackson network as a queue of tasks, of a particular type, which are to be processed by the associated specialized server. A complete pooling of queues, into a single queue, and servers, into a single server, gives rise to an $M/PH/1$ queue, where the server is *flexible* in the sense that it processes all tasks. We assess the value of complete pooling by comparing the steady-state mean sojourn times of these two systems. The main insight from our analysis is that care must be used in pooling. Sometimes pooling helps, sometimes it hurts, and its effect (good or bad) can be unbounded. Also discussed briefly are alternative pooling scenarios, for example complete pooling of only queues which reside in an $M/PH/S$ system, or partial pooling which can be devastating enough to turn a stable Jackson network into an unstable Bransson network. We conclude with some possible future research directions. (*Service Facility Design; Flexible Server; Specialized Server; Service Operations; Efficiency; Stability; Economics of Scale*)

1. Introduction

A fundamental problem in the design and management of stochastic service systems is that of pooling, namely the replacement of several ingredients by a functionally equivalent single ingredient. We analyze the pooling phenomenon within the framework of queueing networks where in our case, as will be explained momentarily, it can take one of three forms: pooling queues (the demand), pooling tasks (the process) or pooling servers (the resources). Here we consider pooling queues and servers simultaneously, but keep the task structure intact, and we provide an efficiency index (5) to determine when such pooling is or is not advantageous.

Our models are described in terms of customers who seek service provided by servers. Service amounts to a collection of tasks, of which there are a finite number of types. Two main models are considered: in the first specialized model, each task type has a server and a queue dedicated to it. For example, Figure 1 exhibits a queueing network in which every customer requires a service that constitutes three tasks, and the tasks are carried out successively, each by its own specialized server. Customers arrive at rate α , average task durations are m_i , and servers' capacities are c_i . In the second

flexible model, servers are capable of handling all tasks and they collectively attend to a single queue of services. For example, Figure 2 exhibits such a model, which arises through pooling the tandem network from Figure 1: customers arrive at rate α , seeking the same three-task service as before; they all join a single queue, which is now attended by a single flexible server of capacity $2c_1 + c_2$.

Customer arrivals are assumed Poisson and task durations exponential. (We comment on these distributional assumptions in the Addendum.) As articulated in §2, we allow a service to consist of a random sequence of tasks in a way that the service duration has a phase-type distribution (a phase corresponds to a task). The specialized (unpooled) model turns out to be a Jackson network (Jackson 1957), as in Figure 3, and the flexible (pooled) architecture is modeled by an $M/PH/1$ system (Neuts 1981) as in Figure 4.

In addition to the above two main models, we also consider briefly alternative designs of pooling. For example, Figure 5 depicts the network from Figure 1, with its queues pooled into a single queue and the servers made flexible while still maintaining their individual identities (see §5.3). Figure 6 depicts partial pooling of

Figure 1
A Specialized Model with Tasks Assigned by Specialized Servers

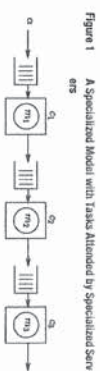


Figure 2
A Flexible Model with Complete Pooling into a Single Queue and a Single Flexible Server

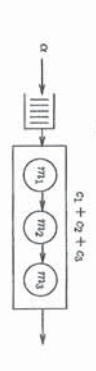


Figure 3
A Specialized Model with Task Repetition and Feedback

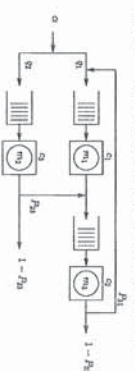
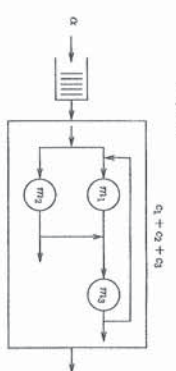


Figure 4
The Flexible Model, under Complete Pooling, that Corresponds to Figure 3



only queues and servers 1 and 2 (see §5.4). Figure 7 depicts a split of the service so that a customer, upon completion of a task, rejoins the queue (see §5.5), and additional designs are possible as well. A common feature of our models is that service is unaltered. For example, in Figures 1, 2, 5, 6, and 7, service always consists of tasks 1, 2, and 3 in succession.

1.1. Motivation

The present research arose from an analysis of a service network consisting of several specialized departments.

MANDELBAUM AND REIMAN
On Pooling in Queuing Networks

The network was redesigned as a pooled single department, which was still responsible for the same services, but whose servers were flexible enough to process all tasks. In trying to analyze this transition, we found that prevalent pooling models failed to cover our network scenario.

Our models provide a new simple framework that helps in assessing the effects on pooling of utilization, variability, and service design. While this is not aimed as a review paper, our framework also relates, as it happens, rather disparate concepts and results, for example (Bransson 1994, Jackson 1957, Klimov 1974, Neuts 1981, Smith and Whitt 1981, and Taha and Pilska 1977). We believe that the usefulness of the framework goes beyond the original motivating applications, pertaining to the design of telephone call centers (Bergandi et al. 1994), evaluation of communication networks (Smith

Figure 5
Complete Pooling of Queues Only (Servers Are Made Flexible but Maintain Individual Identities)

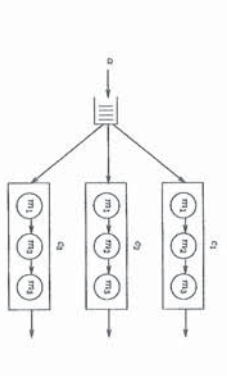


Figure 6
Partial Pooling

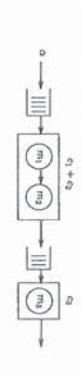
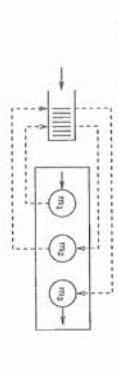


Figure 7
Splitting Services (Each Task Returns to the End of the Queue)



Example: Phase-Type Service Times

Reference: “Length of Stay of Elderly People in Institutional Long-Term Care”, Xie, Chaussalet & Millard, 2005.

Operational significance:

- “Most common causes of delay in **discharge from hospital** are patients awaiting placement in a nursing or residential home and awaiting assessment of their needs.”
- Significant **costs** associated with maintaining elderly people in care homes, hence relevant to “government agencies (funding, planners), insurance companies, and purchasers and providers of care.”

Elderly people go through three states, after being admitted to long-term care:

- **Residential** home care (**R**);
- **Nursing** home care (**N**);
- **Discharge** state (**D**).

Goal: Estimate the **sojourn time in long-term care**, both **duration** and **structure**.

Data: “Paths” of 889 patients, some censored:

- **392 patients**: $R \rightarrow D$ (219 censored);
- **451 patient**: $N \rightarrow D$ (156 censored);
- **46 patients**: $R \rightarrow N \rightarrow D$.

The states **R** and **N** are aggregated: Service time in each is modeled by a *Cozian* (Phase-Type) distribution.

Summary: The above approach is potentially useful in other service contexts. For example, estimating **duration** and **structure** of

- *Telephone or face-to-face services*, in which case data censoring is not important since observations are complete; aggregation is significant, balancing complexity against goodness-of-fit.
- *Customers’ Impatience*, in which case censoring is very important to account for (as will be explained in due time).

A continuous time Markov model for the length of stay of elderly people in institutional long-term care

H. Xie, T. J. Chaussalet and P. H. Millard

University of Westminster, London, UK

[Received January 2003, Final revision January 2004]

Summary. The paper develops a Markov model in continuous time for the length of stay of elderly people moving within and between residential home care and nursing home care. A procedure to determine the structure of the model and to estimate parameters by maximum likelihood is presented. The modelling approach was applied to 4 years' placement data from the social services department of a London borough. The results in this London borough suggest that, for residential home care, a single-exponential distribution with mean 923 days is adequate to provide a good description of the pattern of the length of stay, whereas, for nursing home care, a mixed exponential distribution with means 59 days (short stay) and 794 days (long stay) is required, and that 64% of admissions to nursing home care will become long-stay residents. The implications of these findings and the advantages of the proposed modelling approach in the general context of long-term care are discussed.

Keywords. Length-of-stay modelling; Long-term care; Markov model; Survival

1. Introduction

In the UK, the National Audit Office has recently reported that the most common causes of delay in discharges from hospital are patients awaiting placement in a nursing or residential home and awaiting assessment of their needs (National Audit Office, 2003). Under the 1990 National Health Service and Community Care Act and the Care Standard Act 2000, local authorities in Great Britain are responsible for the placement and finance of adults in publicly funded residential and nursing home care that conforms to national standards. Discharge to long-term care is a central component of plans for acute hospital care and the demand for long-term care is expected to increase substantially as the population ages (Wittenberg *et al.*, 2001). In England, already 1 in 5 people aged 85 years or over live in a long-term care institution (Latho, 2001). In addition, the UK Government is planning to fine local authorities for failing to provide vacancies in residential and nursing home care for hospital discharges. Therefore, it is important for both health authorities and local authorities to have a sound understanding of the patterns of the length of stay (LOS) and movements of residents in long-term care.

A recent survey showed that nearly 70% of the residents in residential and nursing homes were publicly funded and were there permanently (Netten *et al.*, 2001). In earlier research, we found that older people who are placed in nursing homes are more likely to have complex problems. Factors such as being male, immobile, dependent in feeding, urine incontinent, having open wounds and taking multiple drugs are associated with nursing home care placements, whereas older people who are admitted to residential home care are likely to be more independent (Xie

52 H. Xie, T. J. Chaussalet and P. H. Millard
et al., 2002). Therefore, we would expect differences in the pattern of LOS in residential and nursing home care.

Research in the UK shows that the mortality rate for residents in nursing home care is particularly high in the first few months and then gradually levels out (Smith and Lowther, 1976; Bebbington *et al.*, 2001; Roibera *et al.*, 2002). This observation supports the notion of phases in residents' stay in care homes. In the context of hospital geriatric departments, Harrison and Millard (1991) and Taylor *et al.* (1998, 2000) have shown that, despite the great heterogeneity between individuals (Millard, 1988), compartmental and Markov models, which divide patients' LOSs into short-stay and long-stay phases, capture successfully the behaviour of patients' LOSs. Similar results for residential and nursing home care can be expected.

We model the flow of elderly residents within and between residential and nursing home care by using a continuous time Markov model, in which residents' stay in care homes is modelled as a two-phase process: short stay and long stay. First, we describe the model that we propose and present a procedure for determining the model structure and estimating parameters by the method of maximum likelihood. We also show and discuss results that are obtained from fitting the model to a real data set.

2. A model for movement of elderly people in residential and nursing home care

The proposed conceptual model for the movement of elderly people in residential and nursing care facilities is depicted in Fig. 1. In this model, elderly people can be admitted into residential home care or nursing home care directly, either from the community or following discharge from hospital. In each type of care, residents start their stay in the short-stay phase and either leave care after a short period of time or continue their stay to become long-stay residents. People in residential home care can move to nursing home care if their conditions deteriorate to such an extent that residential home care is no longer adequate. In this paper, we consider only those residents who require local authority funding, and we exclude residents whose admissions are meant to provide short respite for their carers. This restriction is imposed because most local authorities have means of determining suitable care placements for applicants requiring public funds; therefore, these admissions will better reflect residents' physical conditions and needs. Movements from nursing home care to residential home care rarely occur among residents who are supported by local authority funds (Bebbington *et al.*, 2001) and are not modelled.

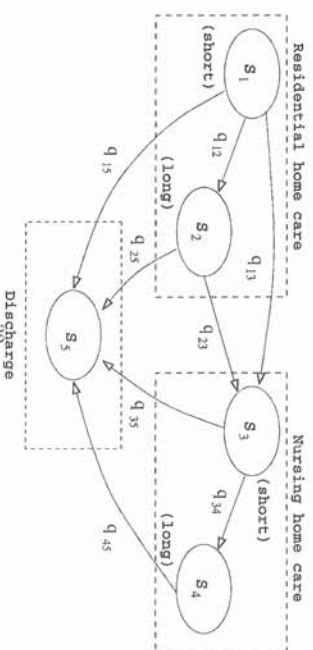


Fig. 1. Markov model for movements of elderly people in residential and nursing home care

Discharges from institutional long-term care are considered permanent. They occur predominantly by death and, although a small number of residents are discharged to the community or hospital, they are not expected to return to institutional long-term care. Discharges to the community are rare for local-authority-funded residents, and those to hospital usually mean terminal care (Bebbington *et al.*, 2001).

We construct a continuous time Markov model of the flow of elderly people within and between residential and nursing home care. The phases in each type of care and the discharges between them form the system states. Given the Markov model that is described in Fig. 1, the generator matrix Q is written as

$$Q = \begin{pmatrix} q_{11} & q_{12} & q_{13} & 0 & q_{15} \\ 0 & q_{22} & q_{23} & 0 & q_{25} \\ 0 & 0 & q_{33} & q_{34} & q_{35} \\ 0 & 0 & 0 & q_{44} & q_{45} \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}, \quad (1)$$

where q_{ij} is the instantaneous transition rate between state i and state j ($i \neq j$), and the elements in the main diagonal are defined such that row sums are 0, i.e. $q_{ii} = -\sum_{j \neq i} q_{ij}$.

3. Maximum likelihood estimation of model parameters

The actual states of the Markov model are not observable. We can only observe which type of care a person is in. For example, at any time, we observe that a person is in residential home care but we do not know whether she or he is in a short-stay (S_1) or long-stay (S_2) state. This is an aggregated Markov process, i.e. a Markov process in which system states are aggregated into a number of classes (Fredkin and Rice, 1986). There are three classes in the model that is outlined in Fig. 1, namely residential home care, nursing home care and discharge (denoted by \mathcal{R} , \mathcal{N} and \mathcal{D} respectively). We partition the matrix \mathbf{Q} according to the class structure of the model, i.e.

$$Q = \begin{pmatrix} Q_{KK} & Q_{KN} & Q_{KD} \\ 0 & Q_{NN} & 0 \\ 0 & 0 & Q_{DD} \end{pmatrix}, \quad (2)$$

where the submatrices correspond to those delimited by broken lines in equation (1) and the subscripts represent system classes. For instance, $Q_{\mathcal{RN}}$ is the submatrix of transition rates from states in \mathcal{R} to states in \mathcal{N} , and $Q_{\mathcal{RR}}$ that of transition rates between states within \mathcal{R} .

The theory or aggregated Markov processes has been motivated by and applied to the modeling of ion channels in neurophysiological applications (Colquhoun and Hawkes, 1981, 1982; Fredkin *et al.*, 1988). Generalization and parameter estimation have been investigated by various researchers, including Ball and Sansom (1989), Fredkin and Rice (1986) and Qin *et al.* (1997). We adapt and modify the approach that was taken by these researchers to suit our modeling needs and to deal with the existence of an absorbing state and censored observations.

3.1. Distribution of sojourn time in a class

Calculating the first-passage time (Cox and Miller, 1965) leads to the probability density function (PDF) of the sojourn time in a class, say class \mathcal{R} (Colquhoun and Hawkes, 1981)

$$f_R(l) = -\phi_R^T \exp^3(Q_{RR}l) Q_{RR}^T l, \quad (3)$$

Table 2. Determination of the number of states in \mathcal{R} and \mathcal{N}

Number of states	Results for residential home care		Results for nursing home care	
	AIC	BIC	AIC	BIC
1	3430.651	3434.733	4879.295	4883.504
2	3433.142	3445.388	4787.788	4787.414
3	3437.142	3457.553	4778.792	4799.835

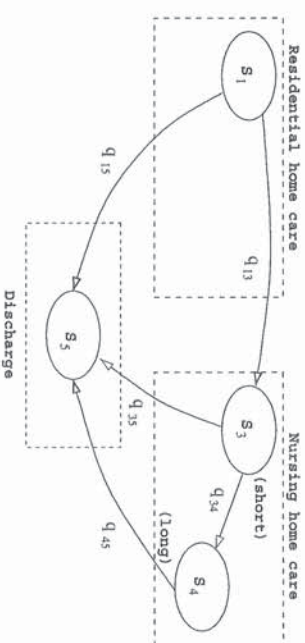


Fig. 2. Structure of the Markov model for the Merton data set

(Fig. 2). The second-stage Markov model fitting procedure converged quickly with the starting-point proposed in Section 3.3. One-dimensional views of the log-likelihood surface along each parameter axis suggested that the maximum was well defined and that the log-likelihood surface was relatively quadratic near the maximum. For each type of care, the close agreement between the survivor curve that was derived from the estimated matrix \hat{Q} (see equation (5)) and the survivor curve that was estimated by the Kaplan–Meier estimator (Kaplan and Meier, 1958) indicates that the Markov model provides a good fit to the data (Fig. 3). This is confirmed by the probability plots (Fig. 4).

4.3. Results

The estimated parameters for the Markov model are summarized in Table 3. These results give interesting insights into the survival patterns of elderly people in institutional long-term care in the London Borough of Merton. A single state provides a good fit to the LOS pattern in residential home care (R), thus indicating a constant rate of departure from R . The average LOS for R is estimated by $1/(q_1 + q_2)$, i.e. 923 days (about 2.5 years). On leaving R , about 79% of the residents will be discharged (permanently) and 21% of them will transfer to nursing home care (N). Two distinctive states are observed in N : a short-stay state with an average LOS of 59 days and a long-stay state with an average LOS of 784 days (about 2.1 years). The rate of discharge from the short-stay state is about five times that from the long-stay state. This agrees with empirical observations that initial mortality is higher for the first few months following admission to nursing care (Smith and Lowthion¹⁹⁷⁶, Bebbington *et al.*, 2001; Rothera *et al.*,

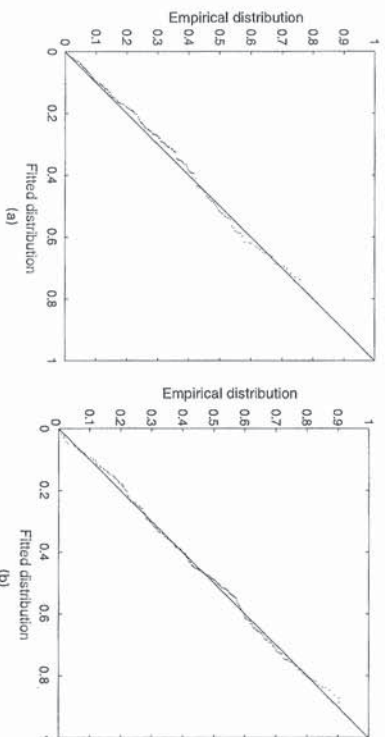


Fig. 4. Probability ($P-P$) plot of the Markov model fitted survivor curves for (a) residential home care and (b) nursing home care for the Merton data set

Table 3. Estimated parameters for the Merton data set

Parameter	Estimate	Standard error	95% confidence interval
q_{13}	0.000228	0.000034	(0.000162, 0.000293)
q_{15}	0.000835	0.000065	(0.000728, 0.000983)
q_{74}	0.010874	0.002961	(0.005071, 0.016677)
q_{75}	0.006138	0.000793	(0.004584, 0.007692)
q_{45}	0.001275	0.000135	(0.001010, 0.001540)

older people who have been placed in \mathcal{R} by the local authority, 50% will stay more than 21 months, 25% will live longer than 3.5 years and 10% will be there after 5.7 years. Of those who have been placed in N , 50% will stay for more than 8 months, 25% will live longer than 2.1 years and 10% will still be there 4.1 years after they have been admitted.

5. Discussion

We have built a continuous time Markov model which captures the flow of elderly people within and between residential and nursing home care. Using the framework of aggregated Markov processes, we derived a procedure for fitting the model to observed data. By modelling the system of long-term care as a whole, we captured the movements between facilities and estimated parameters by using the overall joint likelihood function. Using a real data set we showed that the LOS in residential home care can be approximated by a single-exponential distribution with mean 923 days, whereas in nursing home care a mixed exponential distribution with short-stay mean 59 days and long-stay mean 784 days is needed to provide a good fit. About 21% of residential home care vacancies were created by transfers to nursing home care and 64% of all admissions to nursing home care will become long-stay residents. In nursing home care, the

mortality rate in the short-stay state is about five times that in the long-stay state. Thus, the model quantifies the large heterogeneity in mortality rates that is widely observed in nursing home care.

Extensive research in the UK has been conducted to identify the characteristics that are associated with differences in survival patterns in long-term care. This research has mainly focused on identifying risk factors that are associated with mortality, e.g. Bebbington *et al.* (2001), Dale *et al.* (2001) and Rothera *et al.* (2002). From the point of view of individual elderly people, their doctors and social workers, the identification of risk factors that are associated with transfer, early death and long-term survival is of considerable importance. But, for planning, care managers and budget holders need to know the overall pattern of LOS in long-term care. Our model complements other research in providing a full picture of the overall behaviour of LOS in residential and nursing home care.

Methods that explicitly model the survival time (or the LOS in care) of elderly people have consistently shown that a mixture of exponentials gives a good fit to observed LOS data (Harrison and Millard, 1991; McClean and Millard, 1993; Taylor *et al.*, 1998, 2000). Struthers (1963) first reported that LOS in a hospital geriatric department in Southampton followed a combination of two exponential curves: one had a 'half-life' of 2 months and the other had a half-life of 2 years. A mixed exponential distribution implies that a proportion of elderly people in residential and nursing home care will live substantially longer than the mean and the longer their stay the longer their expected further stay will be. A large proportion of older people who have been placed by the Merton Social Service Department in residential and nursing home care will stay substantially longer than their expected LOS, 2.5 years and 1.5 years respectively. In residential home care, 25% will live longer than 3.5 years and 10% will live longer than 5.7 years. In nursing home care, 25% will live longer than 2.1 years and 10% will live longer than 4.1 years. This means that short-term decisions to increase the number of permanent admissions to residential and nursing home care will have serious long-term financial and organizational consequences. Such action will result in, as time passes, a reduction in the places that are available for new admissions since the number of beds occupied by residents admitted in earlier years increases.

The model that we have developed in this paper could help planning authorities to understand the overall pattern of usage of resources for elderly people in their catchment area. Our model can be extended to cope with possible differences in survival pattern between nursing care residents who are admitted directly and those who are transferred from residential care, although we did not find significant evidence to suggest that such differences existed in the data set that we used. Further work is needed to confirm our findings and to extend the model to take into account the attributes of elderly people, e.g. their age, gender and physical and mental conditions.

Given the importance of having vacancies in long-term care to run acute hospitals efficiently and the significant costs that are associated with maintaining elderly people in care homes, the findings of this paper should be of great interest to Government departments, insurance companies, health and social services planners, and purchasers and providers of residential and nursing home care.

Acknowledgements

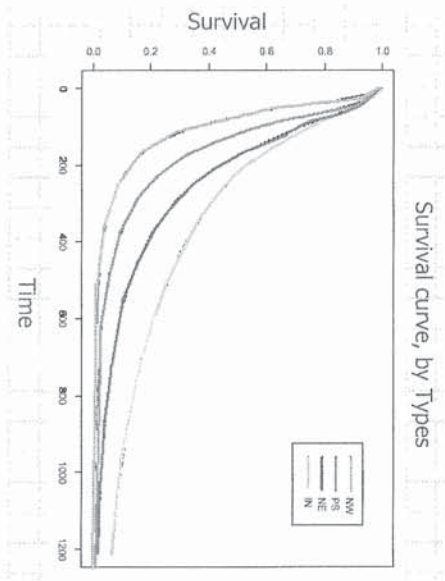
We thank Ms Teresa Temple, Mr Peter Crowther and the late Mr Terry Bucher from the Housing and Social Services Department of the London Borough of Merton for providing the data. This work was partially supported by the Peel Trust and by the Engineering and Physical Sciences Research Council (grant GR/R86430/01).

Comparing Service Durations

First: Means, Standard Deviations, Medians

	Overall	Regular service	New customers	Internet	Stock
Mean	188	181	111	381	269
SD	240	207	154	485	320
Med	114	117	64	196	169

Then: Distributions (Stochastic Order?)



Workload (Offered-Load)

Stationary System Workload (Offered Load):

$$R = \lambda \times E[S]$$

“minutes” of work (=service) that arrive per “minute”.

Example: $\lambda = 3000$ calls/hour; $E[S] = 3$ min.

Consistent time-units, eg. $\lambda = 3000/60 = 50$ calls/min.

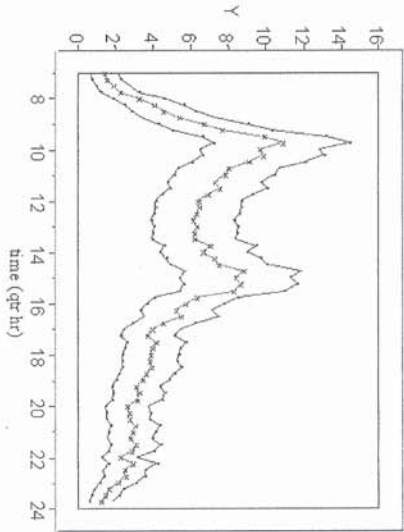
Workload $R = 50 \cdot 3 = 150$ min of work per min.

(If time-units hours? hence Workload in Erlangs.)

Non-Stationary System Workload (Offered Load):

Use $R(t)$ from $M_t/G/\infty$ queue: $R(t) = E[\lambda(t - S_e)] \times E[S]$

Prediction of Workload: Small Israeli Bank



Root Cause Analysis of Emergency Department Crowding and Ambulance Diversion in Massachusetts

A report submitted by the Boston University
Program for the Management of Variability in Health Care Delivery
under a grant from the
Massachusetts Department of Public Health

October, 2002

Emergency Room Diversion Study: Analysis and Findings.

Phase I

Phase I of these investigations involved formulation of a conceptual model that would permit data collection and analysis germane to the problem of ambulance diversion. As preparation for this study, a wide range of relevant medical publications, policy statements and commissioned studies were reviewed. This was followed by personal interviews with representatives in government, hospital administration, public health and the Emergency Medicine community. Information was gathered from throughout Massachusetts and from other key states. Particular attention was given to experience in areas where crowding is particularly severe including metropolitan Boston, San Francisco, Los Angeles and the states of Arizona and Florida. Overall, numerous potential root causes of diversion had been articulated both in the medical literature and lay press, but empirical data to support them were lacking. Available research tended to be descriptive, documenting the extent of crowding without clear delineation of its sources. Various solutions had been proposed and implemented, all without consistent benefit. A partial summary of this analysis has been previously released by the Massachusetts Health Policy Forum of Brandeis University.

An operations management perspective suggested straightforward input-throughput-output analysis. Hospital utilization data provided by the Division of Health Care Finance and Policy was therefore reviewed alongside diversion data provided by regional EMS providers. Analysis of this information revealed the likely operation of mechanisms both internal and external to emergency departments. In addition to simple supply/demand imbalances for emergency care, diversion and utilization patterns suggested bottlenecks and backlogs related to the competition of emergency and non-emergency patients for similar resources. The interrelationships of hospital services then became the focus of attention and patient care pathways were explored with administrators from the two study hospitals.

Two paradigms for the quantitative study of interrelationships among hospital departments were considered. The first involved an analytical approach wherein each relationship was identified, its stochastic character estimated, and appropriate

mathematical models applied. The second involved a simulation approach, wherein stochastic relationships were embedded into computer software that translated real patient flow inputs into utilization and capacity information. Computer simulation was ultimately selected as the route of choice because of its scalability and adaptability.

Phase II

Data Collection/Analysis Effort:

The study was performed at two hospitals in Massachusetts: Hospital A, a large tertiary academic hospital, and Hospital B, a medium-sized acute care community hospital. The following data were collected:

- 42 days of information covering:
 - ‡ 6000+ admissions
 - ‡ 8000+ ED visits
 - ‡ 2000+ staffing/capacity data points
 - ‡ 300,000+ patient movement/status data points

In order to analyze the relationship between diversion status and other factors within the hospital environment all measures were split into observations at one hour increments. The study period of 42 days, with 24 hours per day, yielded a total of 1008 full sets of observations. The analysis required collection of patient flow data well beyond the usual capabilities of contemporary hospital information systems.

Point-biserial coefficients of correlation, with diversion status as the binary variable, were examined against a variety of factors. Comparisons when using full hours of diversion versus partial hours as the "true" condition did not reveal significant differences, so partial diversion hours were evaluated as the "true" binary throughout the analysis for the sake of consistency.

It is important to note that in the real world the decisions to commence or cease diversion status are, but their nature, highly subjective. Because the purpose of the study was to examine the root causes of diversion, we did not approach the task from the standpoint of critiquing or attempting to influence this inherent operational subjectivity. As a result, any such analysis is itself subjective to a certain degree.

Because both hospitals straddled EMS regional borders and diversion rules vary by region, each hospital's data was used for the sake of determining diversion status rather than using centralized EMS data. Also, all diversions were considered equally rather than separately analyzing the factors related to each individual diversion type.

Patterns of diversion were also examined as averages across the hours of the day and the days of the week in order to ascertain relevant hour of the day and day of the week patterns. This data analysis was performed separately for each of the hospitals.

Hospital A:

Diversion Pattern "Hospital A - Diversion Minutes by Hour"

- There were a total of 22 episodes of diversion which started and ended within the study, with an average length of 814 minutes. There was one episode that began prior to the study and ended after the study began and so is not included in this calculation, nor in any calculations which involve the beginning of diversion episodes.
- The hourly diversion pattern shows diversion is highest in the evening hours, settles back down during the early morning hours, and then stays steady until the mid to late afternoon (see Fig. 1).
- The goal of the project was to determine the drivers which create this pattern.

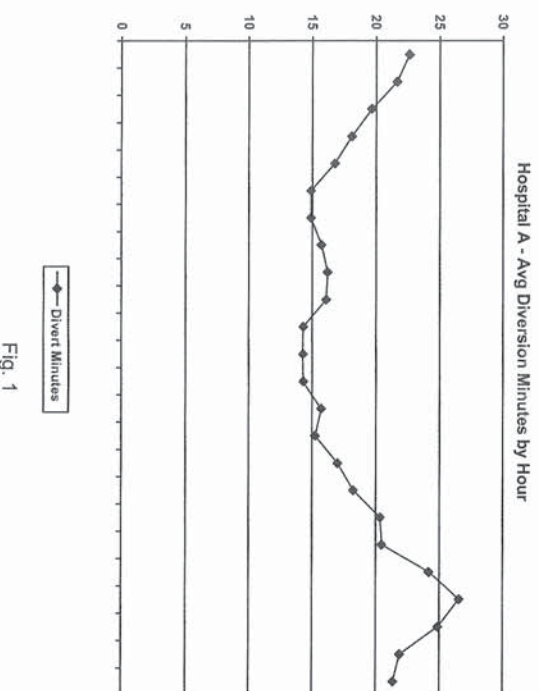


Fig. 1

- The following 3 hypotheses were tested to determine the drivers of diversions:
1. ED arrival rate is too high, leading to diversion when the ED becomes full.
 2. ED processing of patients is too slow, causing backups that lead to diversion
 3. ED arrival and processing rates are fine, but there are not enough beds in the hospital to accommodate the admissions.

There are seven sets of data (see Fig. 2), each representing a different view of arrivals into the ED. The "Arrivals_0" category only includes new arrivals from the hour in question. Each subsequent category, from "Arrivals_1" to "Arrivals_6" includes one more hour's worth added to the total. In other words, "Arrivals_1" includes arrivals from the current hour added to the arrivals from the previous hour, "Arrivals_2" includes all of "Arrivals_1" plus the arrivals from two hours ago, and so on. This is what accounts for the stacked shape as each additional hour is layered on top. Because average length of stay was 340 minutes, 6 hours is used as the maximum lag. Correlation coefficients from each of these cumulatives to Avg Diversion Minutes by hour are as follows:

Arrivals_0 = -0.073
 Arrivals_1 = 0.001
 Arrivals_2 = 0.078
 Arrivals_3 = 0.165
 Arrivals_4 = 0.259
 Arrivals_5 = 0.359
 Arrivals_6 = 0.460

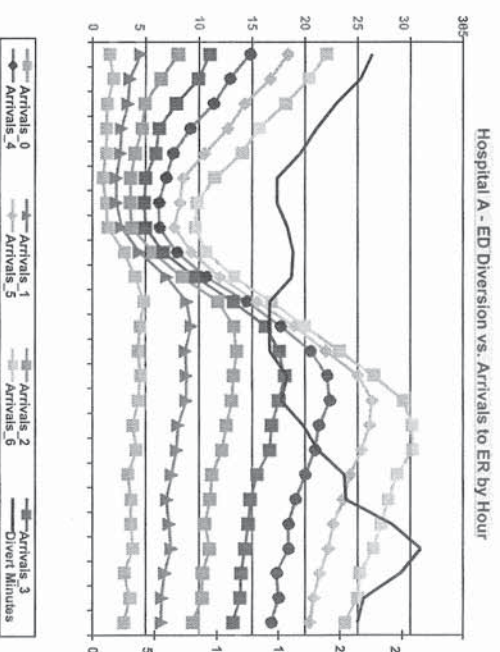


Fig. 2

There is also a possible corollary to hypothesis #1, that overall ED census is a driver of diversion. When counting the non-boarding census and comparing it to diversion status, however, the resulting point-biserial coefficient ($r = -0.051$) makes clear that this potential explanation should be rejected as well.

again points towards examining hospital capacity as the primary determinate of diversion.

Census/Admissions/Discharges: Hospital B.

The overall relationship between inpatient census and ED boarders in Hospital B was similar to that of Hospital A. However, detailed analysis of admission sources in Hospital B is not presented because scheduled demand played a far smaller role than that observed in Hospital A.

During the study period, there were 1,158 weekday unscheduled admissions (average: 38.6/day) and 208 weekday scheduled admissions (average: 6.9/day). This suggests very little operational flexibility in controlling the variability or timing of scheduled arrivals. This likely reflects a fundamental difference between most community hospitals and larger academic centers.

Hospital B Conclusions:

The findings at Hospital B are consistent with and reinforce those at Hospital A. Specifically, there was no evidence that ED process times were temporally or mechanistically related to ED diversion while the relationship between ED arrival rate and diversion was weak. Instead, the data suggest that factors outside of the ED that combine to increase boarders and limit ED capacity are more important.

Phase II Summary:

Detailed flow analysis in two very different types of hospitals yielded similar findings with respect to the root cause of emergency department crowding and ambulance diversion. Neither increased patient inflow nor increased process time could be strongly related to diversion status. Instead, diversion was seen as an outflow problem, with busy emergency departments crowding as patients await transfer to crowded inpatient services. This problem is exacerbated in hospitals with large volumes of scheduled admissions, since these necessarily compete for the same resources. The "collision" of scheduled and unscheduled patient flows results in diversion patterns that are specific and reproducible. Because scheduled patient flows are theoretically controllable, better understanding of this phenomenon may suggest means of decreasing diversion. If the experience here may be generalized, we conclude that institutions with small (or uncontrollable) scheduled patient flows will require addition of resources *on the inpatient side* if diversion is to be substantially reduced.

חלק 2. שימושים.

בבתי הספר ש- $R(t)$, כפי שהוגדר לעיל, משמש כהגדרת העומס המוצע (Offered Load) למערכת שירות בזמן t , כאשר המופע אליה נמשך השירות בה הם כמותו. בחלק 1. (כזכור, $R(t)$ נמדד ב- Erlang).

4.2.1 הסבר בקצרה מדוע אכן ההגדרה מוצגת. עש זאת בהסתמך על הייצוג הרביעי הבא

של $R(t)$:

נסמן $\Lambda(t) -$ מספר המופעים הכולל עד זמן t (כולל).

$\Lambda(t) = E[\Lambda(t)] - E[\Lambda(t)]$ - תחלת מספר המופעים המצטבר עד זמן t .

אזי, מאחת ההגזות של $R(t)$ נובע שמתקיים: $R(t) = E[\Lambda(t) - \Lambda(t - S)]$.

תשובה:

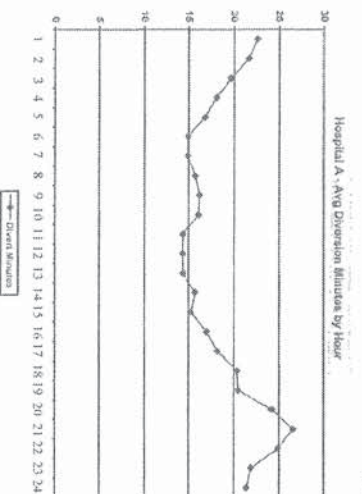
$R(t) =$ כמות העבודה הנמצאת במערכת בזמן t ; עבודה נמדדת ביחידות זמן שירות, או לחלופין

מספר יחידות זמן-שירת הנדרשות לטיפול. לכן $R(t) = E[L(t)]$.

הביטוי $R(t) = E[\Lambda(t) - \Lambda(t - S)]$ מפרש יפה את העבודה שכמות העבודה בזמן t כוללת

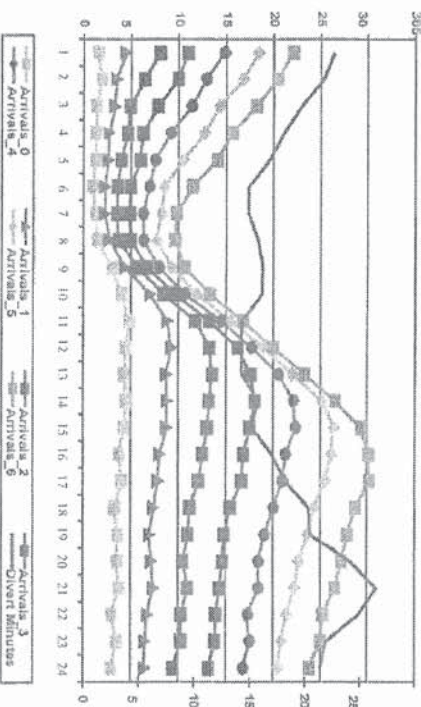
עבודה שהגיעה לפני כן, בהתאם לתחילת ההגעות המצטצ $\Lambda(t)$.

4.2.2 הגרף הבא מוצג את תופעת ה- Ambulance Diversion (הסחת אמבולנסים) שחוצה בכיתה (נופצה בארצות הברית). למשל עד הגרף בשעה 21:00 הוא אחוז הפגמים בתקופה הנמדדת ששעה 20:00-21:00 הוכרח כשעה שבה המיון חסום לכניסת חולים חדשים דרך אמבולנסים (במקרה שלנו 27%).



הגרף הבא מנסה לקשר בין שעת שית הסתת האמבולנסים לשית העומס במיון. הסבר את הגרף. כיצד ניתן היה לפרש את ההתאמה בין הזמנים בהם מתרחשים השיאים? שימו לב: תחלת זמן חשייה במיון תחילת היא 6 שעות.

Hospital A - ED Diversion vs. Arrival to ER by Hour



תשובה:

הגרף $Arrivals_0$ מוצג את כמות החולים הממוצעת לשעה שהגיעו למיון לאורך הימים, עם שית הגעות ב- 11. הגרף $Arrival_1$ בשעה t הוא כמות הגעות בשעה $t +$ כמות הגעות בשעה $t - 1$, כלומר הוא סוכם את הגעות בשעתיים האחרונות. וכי... כלומר הגרף $Arrivals_6$ סוכם את הגעות ב- 6 השעות האחרונות עיי העיסה: $6(t) = \Lambda(t) + \Lambda(t - 1) + \dots + \Lambda(t - 6)$.

לכן גרף 6 ממש בפועל את העומס המוצע $R(t)$ לפי הקירוב הבא:

$6(t) = E[\Lambda(t) - \Lambda(t - S)] \approx \Lambda(t) - E[S] = \Lambda(t) + \Lambda(t - 1) + \dots + \Lambda(t - 6)$ בעזרת

מראים ששית העצמם מתקבל בשעה 17-15, כלומר הוא מוזז בזמן לעומת שית הגעות שקורה ב- 10-11.

אז עדיין יש פער זמן לא מוסבר בין שית העומס הנייל לשית ההסחת (שמסתרש בשעה 12). הפרש זה יכול לגנוי, למשל, מאחד ההסברים הבאים (או משיילבם): א. האמד ל- $R(t)$ לא מספיק מדויק. במקרה זה אנ מכלילים על השיפור הבא: דאז להשתמש בנוסחת העומס המוצע המדויקת, או באחד הקירובים היותר מדויקים שנלמדו בכיתה.

ב. הסתת האמבולנסים, כפי שהוסבר יותר מאוחר, לעומת שעת שית העומס במיון. חסומות, וזה קורה בשעה יותר מאוחר, לעומת שעת שית העומס במיון.

ג. יתכן שפע הזמנים מצביע על כך שהפעלת נוהל הסתת האמבולנסים מופעל מאוחר מדי, כלומר רק לאחר שרואים עומס בולט סביב במיון, מתחילים להפעיל אותו ועד שהוא מגיע לאפסטיכיות מרבית עובר זמן ניכר.