



March 2007

# Operational risk assessment of chemical industries by exploiting accident databases

A. Meel

*University of Pennsylvania*

L. M. O'Neill

*University of Pennsylvania*

J. H. Levin

*University of Pennsylvania*

Warren D. Seider

*University of Pennsylvania, seider@seas.upenn.edu*

U. Oktem

*University of Pennsylvania*

*See next page for additional authors*

Follow this and additional works at: [http://repository.upenn.edu/cbe\\_papers](http://repository.upenn.edu/cbe_papers)

## Recommended Citation

Meel, A., O'Neill, L. M., Levin, J. H., Seider, W. D., Oktem, U., & Karen, N. (2007). Operational risk assessment of chemical industries by exploiting accident databases. Retrieved from [http://repository.upenn.edu/cbe\\_papers/90](http://repository.upenn.edu/cbe_papers/90)

Postprint version. Published in *Journal of Loss Prevention in the Process Industries*, Volume 20, Issue 2, March 2007, pages 113-127. Publisher URL: <http://dx.doi.org/10.1016/j.jlp.2006.10.003>

This paper is posted at ScholarlyCommons. [http://repository.upenn.edu/cbe\\_papers/90](http://repository.upenn.edu/cbe_papers/90)  
For more information, please contact [repository@pobox.upenn.edu](mailto:repository@pobox.upenn.edu).

---

# Operational risk assessment of chemical industries by exploiting accident databases

## Abstract

Accident databases (NRC, RMP, and others) contain records of incidents (e.g., releases and spills) that have occurred in the USA chemical plants during recent years. For various chemical industries, [Kleindorfer, P. R., Belke, J. C., Elliott, M. R., Lee, K., Lowe, R. A., & Feldman, H. I. (2003). Accident epidemiology and the US chemical industry: Accident history and worst-case data from RMP\*Info. *Risk Analysis*, 23(5), 865–881.] summarize the accident frequencies and severities in the RMP\*Info database. Also, [Anand, S., Keren, N., Tretter, M. J., Wang, Y., O'Connor, T. M., & Mannan, M. S. (2006). Harnessing data mining to explore incident databases, the *Journal of Hazardous Material*, 130, 33–41.] use data mining to analyze the NRC database for Harris County, Texas.

Classical statistical approaches are ineffective for low frequency, high consequence events because of their rarity. Given this information limitation, this paper uses Bayesian theory to forecast incident frequencies, their relevant causes, equipment involved, and their consequences, in specific chemical plants. Systematic analyses of the databases also help to avoid future accidents, thereby reducing the risk.

More specifically, this paper presents dynamic analyses of incidents in the NRC database. The NRC database is exploited to model the rate of occurrence of incidents in various chemical and petrochemical companies using Bayesian theory. Probability density distributions are formulated for their causes (e.g., equipment failures, operator errors, etc.), and associated equipment items utilized within a particular industry. Bayesian techniques provide posterior estimates of the cause and equipment-failure probabilities. Cross-validation techniques are used for checking the modeling, validation, and prediction accuracies. Differences in the plant- and chemical-specific predictions with the overall predictions are demonstrated. Furthermore, extreme value theory is used for consequence modeling of rare events by formulating distributions for events over a threshold value. Finally, the fast-Fourier transform is used to estimate the capital at risk within an industry utilizing the *frequency* and *loss-severity* distributions.

## Keywords

risk, frequency modeling, consequence modeling, abnormal events, chemical plants

## Comments

Postprint version. Published in *Journal of Loss Prevention in the Process Industries*, Volume 20, Issue 2, March 2007, pages 113-127. Publisher URL: <http://dx.doi.org/10.1016/j.jlp.2006.10.003>

## Author(s)

A. Meel, L. M. O'Neill, J. H. Levin, Warren D. Seider, U. Oktem, and N. Karen



# Operational risk assessment of chemical industries by exploiting accident databases

A. Meel<sup>a</sup>, L.M. O'Neill<sup>a</sup>, J.H. Levin<sup>a</sup>, W.D. Seider<sup>a,\*</sup>, U. Oktem<sup>b</sup>, N. Keren<sup>c</sup>

<sup>a</sup>Department of Chemical and Biomolecular Engineering, University of Pennsylvania, Philadelphia, PA 19104-6393, USA

<sup>b</sup>Risk Management and Decision Center, Wharton School, University of Pennsylvania, Philadelphia, PA 19104-6340, USA

<sup>c</sup>Department of Agricultural and Biosystems Engineering, Iowa State University, Ames, IA 50011-3080, USA

Received 11 June 2006; received in revised form 17 October 2006; accepted 18 October 2006

## Abstract

Accident databases (NRC, RMP, and others) contain records of incidents (e.g., releases and spills) that have occurred in the USA chemical plants during recent years. For various chemical industries, [Kleindorfer, P. R., Belke, J. C., Elliott, M. R., Lee, K., Lowe, R. A., & Feldman, H. I. (2003). Accident epidemiology and the US chemical industry: Accident history and worst-case data from RMP\*Info. *Risk Analysis*, 23(5), 865–881.] summarize the accident frequencies and severities in the RMP\*Info database. Also, [Anand, S., Keren, N., Tretter, M. J., Wang, Y., O'Connor, T. M., & Mannan, M. S. (2006). Harnessing data mining to explore incident databases. *Journal of Hazardous Material*, 130, 33–41.] use data mining to analyze the NRC database for Harris County, Texas.

Classical statistical approaches are ineffective for low frequency, high consequence events because of their rarity. Given this information limitation, this paper uses Bayesian theory to forecast incident frequencies, their relevant causes, equipment involved, and their consequences, in specific chemical plants. Systematic analyses of the databases also help to avoid future accidents, thereby reducing the risk.

More specifically, this paper presents dynamic analyses of incidents in the NRC database. The NRC database is exploited to model the rate of occurrence of incidents in various chemical and petrochemical companies using Bayesian theory. Probability density distributions are formulated for their causes (e.g., equipment failures, operator errors, etc.), and associated equipment items utilized within a particular industry. Bayesian techniques provide posterior estimates of the cause and equipment-failure probabilities. Cross-validation techniques are used for checking the modeling, validation, and prediction accuracies. Differences in the plant- and chemical-specific predictions with the overall predictions are demonstrated. Furthermore, extreme value theory is used for consequence modeling of rare events by formulating distributions for events over a threshold value. Finally, the fast-Fourier transform is used to estimate the capital at risk within an industry utilizing the *frequency* and *loss-severity* distributions.

© 2006 Published by Elsevier Ltd.

**Keywords:** Risk; Frequency modeling; Consequence modeling; Abnormal events; Chemical plants

**Abbreviations:** Companies A, B, C, D, E, F, G, A, B, C, D, E, F, G; Basic indicator approach, BIA; Capital at risk, CaR; Center for chemical process safety (AIChE), CCPS; Equipment failure, EF; Environmental protection agency, EPA; Extreme value theory, EVT; Fast-Fourier transform, FFT; Heat transfer units, HT; Inverse fast-Fourier transform, IFFT; Internal measurement approach, IMA; Loss distribution approach, LDA; Markov-chain Monte Carlo, MCMC; Major accident reporting system, MARS; National response center, NRC; Others, O; Operator error, OE; Occupational safety and health administration, OSHA; Process safety incident database, PSID; Process safety management, PSM; Process units, PU; Process vessels, PV; Quantile-quantile, Q-Q; Risk management plan, RMP; Standardized approach, SA; Storage vessel, SV; Transfer line, TL

\*Corresponding author. Tel.: +1 215 898 7953.

E-mail address: seider@seas.upenn.edu (W.D. Seider).

0950-4230/\$ - see front matter © 2006 Published by Elsevier Ltd.  
doi:10.1016/j.jlpp.2006.10.003

## 1. Introduction

Since the accidents at Flixborough, Seveso, and Bhopal, the reporting of abnormal events in the chemical industries has been encouraged to collect accident precursors. Efforts to increase the reporting of near-misses, with near-miss management audits, have been initiated by the Wharton Risk Management Center (Phimister, Oktem, Kleindorfer, & Kunreuther, 2003). In addition, the AIChE center for chemical process safety (CCPS) has facilitated the development of a process safety incident database (PSID) to collect and share incident information, permitting indus-

## Nomenclature

	$N_{\text{total}}$	total number of incidents	59
	$N_{\text{U}}$	number of incidents associated with unknown causes	61
3	$a, b$	parameters of <i>Beta</i> prior probability distribution	
5	$a_i, b_i$	parameters of prior probability distribution of cause $i$ for an incident	63
7	$d_1, d_2, d_3$	cumulative number of incidents of causes EF, OE, and O at the end of each year	65
9	$e_i$	probability of involvement of equipment type $i$	67
11	$E(\mu Data)$	expected posterior mean of $\mu$	
13	$E(q Data)$	expected posterior mean of $q$	69
15	$E(y)$	expected value of number of incidents in a year	
17	$E[Y_i Y_{-i}]$	expected value of prediction of incident in year $i$ based on incidents in $Y_{-i}$	71
19	$f(e_i)$	prior probability distribution of involvement of equipment $i$ for an incident	73
21	$f(x_i Data)$	posterior probability distribution of involvement of equipment $i$ conditional upon <i>Data</i>	75
23	$f(x_i)$	prior probability distribution of cause $i$ for an incident	
25	$f(x_i Data)$	posterior probability distribution of cause $i$ conditional upon <i>Data</i>	
27	$f_l$	discrete <i>loss-severity</i> distribution function	
29	$f_z(Z)$	discrete probability distribution function of total loss	
31	$F_u(y)$	cumulative probability distribution for distribution of losses, $l$ , over threshold $u$	81
33	$G(l)$	<i>Generalized Pareto</i> distribution of losses	
35	$l$	loss associated with an incident	
37	$M_i + N_i + O_i$	cumulative number of incidents associated with equipment $i$ at the end of each year	87
39	$n_p$	number of points desired in <i>total loss</i> distribution	
41	$N_{C/P}$	number of incidents associated with compressors and pumps	
43	$N_d$	amount of damage, \$	
45	$N_e$	number of evacuations	
47	$N_{EF}$	number of incidents associated with equipment failures	
49	$N_f$	number of fatalities	
51	$N_h$	number of hospitalizations	
53	$N_{HT}$	number of incidents associated with heat-transfer equipment items	
55	$N_i$	number of injuries	
57	$N_{OE}$	number of incidents associated with operator errors	
	$N_{PU}$	number of incidents associated with process units	
	$N_{SV}$	number of incidents associated with storage vessels	
	$N_t$	number of years	
	$N_{TL}$	number of incidents associated with transfer-line equipment	
	$p(\lambda)$	prior distribution of $\lambda$	
	$p(\lambda Data)$	posterior distribution of $\lambda$ given <i>Data</i>	
	$p(q Data)$	marginal posterior distribution of $q$ given <i>Data</i>	
	$p(\mu Data)$	marginal posterior distribution of $\mu$ given <i>Data</i>	
	$P_N$	probability generating function of the frequency of events, $N$	
	$p_i, q_i$	parameters of prior probability distribution of involvement of equipment $i$ in an incident	
	$q$	parameter of the <i>Negative Binomial</i> distribution	
	$s$	total number of incidents in $N_t$ years	
	$u$	threshold value of $l$ for <i>loss-severity</i> distribution	
	$V(y)$	variance of number of incidents per year	
	$w_d$	dimensionless damage measure	
	$w_e$	dollar amount per evacuation, \$	
	$w_f$	dollar amount per fatality, \$	
	$w_h$	dollar amount per hospitalization, \$	
	$w_i$	dollar amount per injury, \$	
	$x_1, x_2, x_3$	probabilities of causes EF, OE, and O for an incident	
	$y_i$	number of incidents in year $i$	
	$z_i$	predictive score for incidents in year $i$	
	$Z$	total annual loss for a company	
	<i>Greek</i>		
	$\alpha, \beta$	parameters for <i>Gamma</i> density distribution function	
	$\beta(a, b)$	<i>Beta</i> density distribution with parameters $a$ and $b$	
	$\phi_l$	characteristic function of the <i>loss-severity</i> distribution	
	$\phi_Z$	characteristic function of <i>total loss</i> distribution	
	$\lambda$	average annual number of incidents	
	$\lambda_B$	average annual number of incidents for company B with losses greater than $u$	
	$\lambda_F$	average annual number of incidents for company F with losses greater than $u$	
	$\mu$	parameter of the <i>Negative Binomial</i> distribution	
	$\xi, \beta$	parameters of the <i>generalized Pareto</i> distribution	
	$\gamma(\alpha, \beta)$	<i>Gamma</i> distribution with parameters $\alpha$ and $\beta$	
	<i>Subscript</i>		
	$i$	year counter	
	$n$	year vector	

1 trial participants access to the database, while sharing their  
 2 collective experiences (CCPS, 1995). Finally, the Mary Kay  
 3 Safety Center at Texas A&M University (TAMU) has been  
 4 gathering incident data in the chemical industries (Anand  
 5 et al., 2006; Mannan, O'Connor, & West, 1999).

6 An incident database, involving oil, chemical, and  
 7 biological discharges into the environment in the USA  
 8 and its territories, is maintained by the national response  
 9 center (NRC) (NRC, 1990). While companies participate  
 10 voluntarily, raising reliability concerns, the NRC database  
 11 for Harris County, Texas, is acknowledged to be reliable  
 12 thanks to the conscientious efforts of many chemical  
 13 companies in reporting incidents. Moreover, the Mary Kay  
 14 Safety Center has concentrated time and resources toward  
 15 refining the Harris County database to increase its  
 16 reliability and consistency.

17 To record accidents, European industries submit their  
 18 data to the major accident-reporting system (MARS)  
 19 (Rasmussen, 1996), while a database for chemical compa-  
 20 nies in the USA is created from risk management plans  
 21 (RMP) submitted by facilities subject to Environmental  
 22 protection agency's (EPA) chemical accidental release  
 23 prevention and response regulations (Kleindorfer et al.,  
 24 2003; RMP, 2000).

25 Several researchers have been analyzing and investigat-  
 26 ing incident databases to identify common trends and to  
 27 estimate risks. For example, Chung and Jefferson (1998)  
 28 have developed an approach to integrate accident data-  
 29 bases with computer tools used by chemical plant  
 30 designers, operators, and maintenance engineers, permit-  
 31 ting accident reports to be easily accessed and analyzed. In  
 32 addition, Sonnemans, Korvers, Brombacher, van Beek,  
 33 and Reinders (2003) have investigated 17 accidents that  
 34 have occurred in the Netherlands petrochemical industries  
 35 and have demonstrated qualitatively that had accident  
 36 precursor information been recorded, with proper mea-  
 37 sures to control future occurrences, these accidents could  
 38 have been foreseen and possibly prevented. Furthermore,  
 39 Sonnemans and Korvers (2006) observed that even after  
 40 recognizing accident precursors and disruptions, the  
 41 operating systems inside companies often fail to prevent  
 42 accidents. The results of yet another analysis feature the  
 43 lessons learned from the major accident and near-miss  
 44 events in Germany from 1993 to 1996 (Uth, 1999; Uth &  
 45 Wiese, 2004). Finally, Elliott, Wang, Lowe, and Kleindor-  
 46 fer (2004) analyzed the frequency and severity of accidents  
 47 in the RMP database with respect to socioeconomic factors  
 48 and found that larger chemically intensive companies are  
 49 located in counties with larger African-American popula-  
 50 tions and with both higher median incomes and higher  
 51 levels of income inequality. Note that accident precursors  
 52 have been studied also in railways, nuclear plants, health  
 53 science centers, aviation, finance companies, and banking  
 54 systems.

55 On the risk estimation frontier, Kirchsteiger (1997)  
 56 discussed the strengths and weaknesses of probabilistic  
 57 and deterministic methods in risk analysis using illustra-

58 tions associated with nuclear and chemical plants. It is  
 59 argued that probabilistic methods are more cost-effective,  
 60 giving results that are easier to communicate to decision  
 61 and policy makers. In addition, Goossens and Cooke  
 62 (1997) described the application of two risk assessment  
 63 techniques involving: (i) formal expert judgment to  
 64 establish quantitative subjective assessments of design  
 65 and model parameters, and (ii) system failure analysis,  
 66 with accident precursors, using operational evidence of  
 67 system failures to derive the failure probability of the  
 68 system. Furthermore, a human and organizational reli-  
 69 ability analysis in accident management (HORAAM)  
 70 method was introduced to quantify human and organiza-  
 71 tional factors in accident management using decision trees  
 72 (Baumont, Menage, Schneiter, Spurgin, & Vogel, 2000).

73 In this work, statistical methods are introduced to  
 74 estimate the operational risk for seven companies, includ-  
 75 ing petrochemical and specialty chemical manufacturers,  
 76 using the NRC database for Harris County, with the risk  
 77 estimated as the product of the frequency and conse-  
 78 quences of the incidents. Fig. 1 shows the algorithm for  
 79 calculating the operational risk of a chemical company.  
 80 For a company in the database, the incidents are extracted  
 81 on a yearly basis. Then, the frequency distribution of the  
 82 incidents is estimated using a  $\gamma$ -Poisson Bayesian model.  
 83 Note that significant differences in the prediction of  
 84 incidents are observed for the individual companies, as  
 85 compared with predictions obtained when the incidents  
 86 from all of the companies are lumped together. The  
 87 Bayesian theory upgrades prior information available, if  
 88 any, using data to increase the confidence level in modeling  
 89 the frequency of incidents, decreasing the uncertainty in  
 90 decision-making with annual information upgrades (Ro-  
 91 bert, 2001).

92 Additional  $\gamma$ -Poisson Bayesian models are developed to  
 93 provide the frequency distribution of the day of the week  
 94 on which the incidents occur, the equipment types  
 95 involved, the causes behind the incidents, and the chemicals  
 96 involved. In parallel, the failure probabilities of the process  
 97 units, as well as the causes of the incidents, are predicted  
 98 using a  $\beta$ -Bernoulli Bayesian model.

99 Later, a *loss-severity* distribution of the incidents is  
 100 modeled using extreme value theory (EVT) by formulating  
 101 a quantitative index for the loss as a weighted sum of the  
 102 different types of consequences. Through EVT, both  
 103 extreme and unusually rare events, which characterize  
 104 incidents reported in the chemical industries, are modeled  
 105 effectively. Note that EVT has been applied in structural,  
 106 aerospace, ocean, and hydraulic engineering (Embrechts,  
 107 Kluppelberg, & Mikosch, 1997). Herein, EVT is introduced  
 108 to measure the operational risk in the chemical industries.

109 Finally, the operational risk of the individual chemical  
 110 industries is computed by performing fast-Fourier trans-  
 111 forms (FFT) of the product of the *frequency* and *loss-*  
 112 *severity* distributions to obtain the *total loss* distribution  
 113 and the capital at risk (CaR). This approach to measuring  
 risks in specific companies provides a quantitative frame-



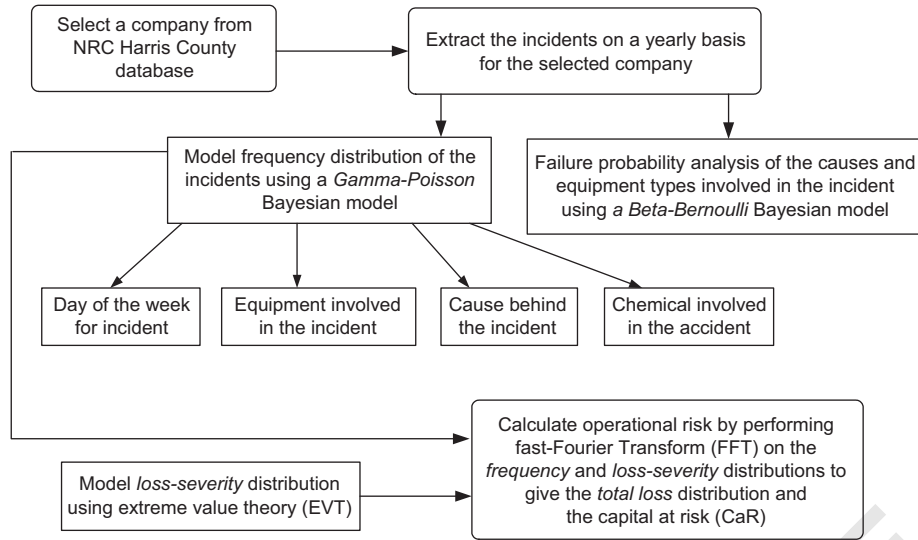


Fig. 1. Algorithm to calculate the operational risk of a chemical company.

work for decision-making at higher levels. Using the platform provided, the chemical industries should be encouraged to collect accident precursor data more regularly. Through implementation of this dynamic risk assessment methodology, improved risk management strategies should result. Also, the handling of third party investigations should be simplified after accidents.

To begin the detailed presentation of this algorithm, Section 2 describes the concepts of Bayesian theory for prediction of the numbers of incidents annually. Then, the NRC database, the Bayesian predictive models, and the *loss-severity* distribution using EVT, are described in Section 3. The CaR calculations using FFTs are discussed in Section 4. Finally, conclusions are presented in Section 5.

## 2. Modeling the frequency of incidents

Bayesian theory is helpful in formulating the annual frequency of occurrence of incidents for a company. The relationship between the mean and the variance of the annual incidents, over many years, determines the best choice of distribution. For example, the *Poisson* distribution is suitable when the mean and the variance of the data are in close proximity. When the predictions of the *Poisson* distribution are poor, other distributions are used; for instance, the *Negative Binomial* distribution, when the variance exceeds the mean (Bradlow, Hardie, & Fader, 2002).

### 2.1. Poisson distribution

The annual number of occurrences of an incident is a non-negative, integer-valued outcome that can be estimated using the *Poisson* distribution for  $y$ :

$$y \sim p(y = y_i) = \left\{ \frac{\lambda^{y_i} e^{-\lambda}}{y_i!} \right\}, \quad y_i \in \{I^1\}, y_i \geq 0, \quad \lambda > 0, \quad (1a)$$

where  $y_i$  is the number of incidents in year  $i$ , and  $\lambda$  is the annual average number of incidents, with the expected value,  $E(y)$ , and variance,  $V(y)$ , equal to  $\lambda$ . Due to uncertainty, the prior distribution for  $\lambda$  is assumed to follow a  $\gamma$ -distribution,  $\lambda \sim \gamma(\alpha, \beta)$ :

$$p(\lambda) \propto \lambda^{\alpha-1} e^{-\beta\lambda}, \quad \alpha > 0, \quad \beta > 0. \quad (1b)$$

From Baye's theorem, the posterior distribution,  $p(\lambda|Data)$ , is:

$$p(\lambda|Data) \propto l(Data|\lambda)p(\lambda) \propto (\lambda^s e^{-N_t\lambda}) \times (\lambda^{\alpha-1} e^{-\beta\lambda}) \propto \lambda^{(\alpha+s)-1} e^{-(\beta+N_t)\lambda}, \quad (1c)$$

where  $Data = (y_0, y_1, \dots, y_{N_t})$ ,  $s = \sum_{i=0}^{N_t} y_i$ ,  $N_t$  is the number of years, and  $l(Data|\lambda)$  is the *Poisson* likelihood distribution. Note that  $p(\lambda|Data)$  is also a *Gamma* distribution,  $\gamma(\alpha+s, \beta+N_t)$ , because  $\lambda$  is distributed according to  $\gamma(\alpha, \beta)$ , which is a conjugate prior to the *Poisson* distribution. The mean of the posterior distribution is the weighted average of the means of the prior and likelihood distributions:

$$\frac{\alpha + s}{\beta + N_t} = \frac{\beta}{\beta + N_t} \left( \frac{\alpha}{\beta} \right) + \frac{N_t}{\beta + N_t} \frac{s}{N_t}, \quad (1d)$$

and the variance of the posterior distribution is  $(\alpha + s)/(\beta + N_t)^2$ .

The predictive distribution to estimate the number of incidents in the next year,  $y_{N_t+1}$ , conditional on the observed  $Data$ , is discussed by Meel and Seider (2006). This gives a predictive mean,  $(\alpha + s)/(\beta + N_t)$ , and predictive variance,  $(\alpha + s)/(\beta + N_t)[1 + 1/(\beta + N_t)]$ , and consequently, the posterior and predictive means are the same, while the predictive variance exceeds the posterior variance.

## 2.2. Negative binomial distribution

The annual number of occurrences of an incident is a non-negative, integer-valued outcome that can be estimated using the *Negative Binomial* distribution for  $y$ :

$$y \sim (q)^\mu (1-q)^{y_i} \quad y_i \in \{I^1\}, y_i \geq 0, \quad \mu > 0, \quad q \geq 0, \quad (1e)$$

where  $y_i$  is the number of incidents in year  $i$ , and  $\mu(1-q)/q$  is the expected annual (mean) number of incidents,  $E(y)$ , and  $\mu(1-q)/q^2$  is the expected variance,  $V(y)$ . Due to uncertainty, the prior distribution for  $\mu$  is assumed to follow a *Gamma* distribution,  $\mu \sim \gamma(\alpha, \beta)$ :

$$p(\mu) \propto \mu^{\alpha-1} e^{-\beta\mu}, \quad \alpha > 0, \quad \beta > 0, \quad (1f)$$

and that for  $q$  is assumed to follow a *Beta* distribution,  $q \sim \beta(a, b)$ :

$$p(q) \propto q^{a-1} (1-q)^{b-1}, \quad a > 0, \quad b > 0. \quad (1g)$$

From Baye's theorem, the posterior distribution,  $p(\mu, q | Data)$ , is

$$\begin{aligned} p(\mu, q | Data) &\propto l(Data | \mu, q) p(\mu) p(q) \\ &\propto q^{n\mu} (1-q)^s (\mu^{\alpha-1} e^{-\beta\mu}) q^{a-1} (1-q)^{b-1} \\ &\propto q^{n\mu+a-1} (1-q)^{s+b-1} (\mu^{\alpha-1} e^{-\beta\mu}), \end{aligned} \quad (1h)$$

where  $Data = (y_0, y_1, \dots, y_{N_t})$ ,  $s = \sum_{i=0}^{N_t} y_i$ ,  $N_t$  is the number of years, and  $l(Data | \mu, q)$  is the *Negative Binomial* likelihood distribution. The marginal posterior distributions,  $p(\mu | Data)$  and  $p(q | Data)$ , and the posterior means  $E(\mu | Data)$  and  $E(q | Data)$  are obtained using the Markov Chain Monte-Carlo (MCMC) method in the WINBUGS software (Spiegelhalter et al., 2003). These added calculations are not needed for the *Poisson* distribution, in which the expected value,  $E(\lambda | Data)$ , is computed easily using Eq. (1d).

## 2.3. Model-checking

To check the accuracy of the model, the number of incidents in year  $i$ ,  $y_i$ , is removed, leaving the data,  $y_{-i} = (y_0, \dots, y_{i-1}, y_{i+1}, \dots, y_{N_t})$ , over  $N_t - 1$  years. Then, a Bayesian model applied to  $y_{-i}$  is used to predict  $y_i$ . Finally,  $y_i$  and  $E[y_i | y_{-i}]$  are compared, and predictive  $z$ -scores are used to measure their proximity:

$$z_i = \frac{y_i - E[y_i | y_{-i}]}{\sqrt{V[y_i | y_{-i}]}}. \quad (2)$$

For a good model, the mean and standard deviation of  $z = (z_0, \dots, z_{N_t})$  should approach zero and one, respectively.

## 3. Analysis of the NRC database

The NRC database contains reports on the oil, chemical, radiological, biological, and etiological discharges into the environment in the USA and its territories (NRC, 1990). A typical incident report includes the date of the incident, the

chemical involved, the cause of the incident, the equipment involved, the volume of the chemical release, and the extent of the consequences. Herein, the incidents reported for Harris County, Texas, for seven specific facilities during the years 1990–2002, are analyzed to determine their frequencies and consequences (loss-severities). This dataset was obtained from the Mary Kay Safety Center at TAMU, which filtered the NRC database for Harris County, taking care to eliminate duplications of incidents when they occurred. More specifically, the filtered dataset by Anand et al. (2006), comprised of 7265 records, is used for further processing.

The equipment is classified into 13 major categories: electrical equipment ( $E_1$ ), pumps/compressors ( $E_2$ ), flare stacks ( $E_3$ ), heat-transfer equipment ( $E_4$ ), hoses (flexible pipes) ( $E_5$ ), process units ( $E_6$ ), process vessels (PV) ( $E_7$ ), separation equipment ( $E_8$ ), storage vessels ( $E_9$ ), pipes and fittings ( $E_{10}$ ), unclassified equipment ( $E_{11}$ ), relief equipment ( $E_{12}$ ), and unknowns ( $E_{13}$ ). The Harris County database includes several causes of the incidents, including equipment failures (EF), operator errors (OE), unknown causes (U), dumping (intentional and illegal deposition of material on the ground), and others, with the EF and OE causes being the most significant. Herein, the unknown causes (U), dumping, and others are combined and referred to as others (O).

### 3.1. Prediction of incidents at chemical companies

Table 1 shows the number of incidents extracted from the NRC database for the seven companies. The total number of incidents,  $N_{total}$ , and the number of incidents of EF,  $N_{EF}$ , OE,  $N_{OE}$ , and due to unknown causes,  $N_U$ , are listed during the years 1990–2002. In addition, from the 13 equipment categories, the number of incidents of process units,  $N_{PU}$ , storage vessels,  $N_{SV}$ , compressors/pumps,  $N_{C/P}$ , heat-transfer equipment,  $N_{HT}$ , and transfer-line equipment,  $N_{TL}$ , are included. Note that the large excess of EF compared with the numbers of OE was unanticipated. Perhaps this is due to cost-saving measures that have reduced maintenance budgets, with major repairs postponed until they are deemed to be urgent. Also, because automated equipment often experiences fewer failures than those related to the inconsistencies of the operators, it is likely that many reported EF are indirectly a result of OE.

For each of the seven companies, the numbers of incidents were predicted for future years utilizing data from previous years. Included are the total number of incidents,  $N_{total}$ , the number of incidents associated with each equipment type, and the number of incidents associated with each cause. In the remainder of this section, selected results are presented and discussed.

Figs. 2(a) and (b) show the predictions of the number of incidents for companies B and F using *Poisson* distributions which are chosen arbitrarily to illustrate the variations in the predictive power of the models. In these figures, the number of incidents for the year  $n$  are forecasted using

1 Table 1  
 Number of incidents for seven companies in the NRC database

3 Companies	Type	$N_{total}$	$N_{EF}$	$N_{OE}$	$N_U$	$N_{PU}$	$N_{SV}$	$N_{C/P}$	$N_{HT}$	$N_{TL}$
5 A	Petrochemical	688	443	56	101	59	101	86	58	121
B	Petrochemical	568	387	48	88	110	69	127	47	56
7 C	Specialty chemical	401	281	35	46	45	61	10	28	77
D	Petrochemical	220	122	24	16	25	16	36	27	15
E	Specialty chemical	119	77	21	8	13	22	11	12	23
9 F	Specialty chemical	83	57	14	7	6	21	8	10	18
G	Specialty chemical	18	9	2	5	1	1	1	3	2

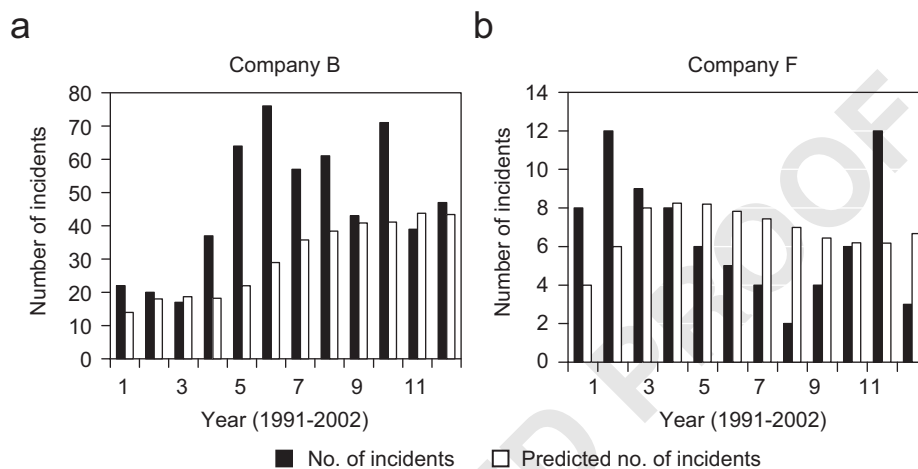


Fig. 2. Total number of incidents: (a) company B, (b) company F.

the *Gamma-Poisson* Bayesian techniques based on the number of incidents from 1990 to  $n-1$ , where  $n = 1991, 1992, \dots, 2002$ . These are compared to the number of incidents that occurred in year  $n$  for companies B and F, respectively.

In the absence of information to model the prior distribution for the year 1990,  $\alpha$  and  $\beta$  are assumed to be 0.001, providing a relatively flat distribution in the region of interest; that is, a non-informative prior distribution. Note that information upon which to base the prior parameters would enhance the early predictions of the models. This has been illustrated for a *Beta-Bernoulli* Bayesian model, using informative and non-informative prior distributions, showing the sensitivity of the predictions to the prior values (Meel & Seider, 2006). For company B, using non-informative prior distributions, either the numbers of incidents are close to the predicted numbers or higher than those predicted. However, for company F, the numbers of incidents are close to or less than those predicted.

When examining the results for the seven companies, the sizable variations in the number of incidents observed in a particular year are attributed to several factors including management and planning efforts to control the incidents, it being assumed that no significant differences occurred to affect the reporting of the incidents from 1990 to 2002—

although OSHA's PSM standard and EPA's RMP rule were introduced in 1992 and 1996, respectively. Therefore, when the number of incidents is less than those predicted, it seems clear that good incident-control strategies were implemented within the company. Similarly, when the number of incidents is higher than those predicted, the precursor data yields a warning to consider enhancing the measures to reduce the number of incidents in the future.

A good agreement between the numbers of incidents predicted and observed indicates that a *stable equilibrium* is achieved with respect to the predictive power of the model. Such a state is achieved when the numbers of incidents and their causes do not change significantly from year-to-year. Note, however, that even as stable equilibrium is approached, efforts to reduce the number of incidents should continue. This is because, even when successful measures are taken year after year (that reduce the number of incidents), the predictive values are usually conservative, lagging behind until the incidence rates converge over a few years.

Next, the results of the Bayesian model checking using the *R* software package (Gentleman et al., 2005) to compute predictive distributions are presented in quantile-quantile ( $Q-Q$ ) plots. For company F, Fig. 3(a) shows the density profile of incidents, while Fig. 3(b) shows the normal  $Q-Q$  plot, which compares the distribution of  $z$



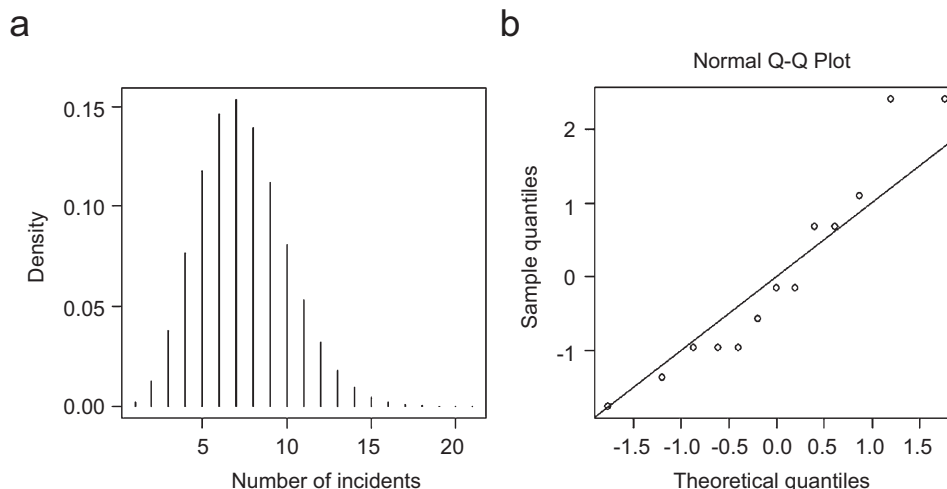


Fig. 3. Company F: (a) density of incidents, (b)  $Q-Q$  plot.

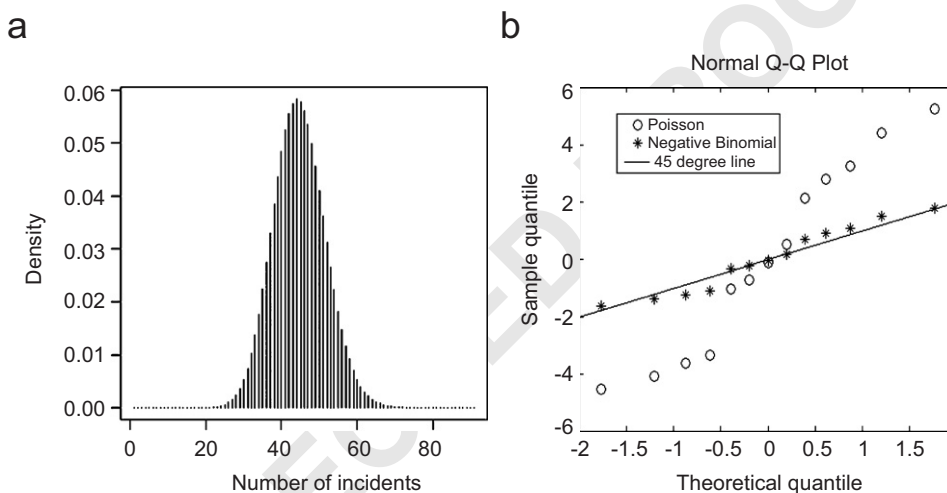


Fig. 4. Company B: (a) density of incidents, (b)  $Q-Q$  plot.

(Eq. (2)) to the normal distribution (represented by the straight line), where the elements of  $z$  are represented by circles. The sample quantiles of  $z$  (ordered values of  $z$ , where the elements,  $z_{i_s}$  are called quantiles) are close to the theoretical quantiles (equally-spaced data from a normal distribution), confirming the accuracy of the model predictions. Most of the values are in good agreement, except for two outliers at the theoretical quantiles, 1.0 and 1.5.

Figs. 4(a) and (b) show the density profile of incidents and the  $Q-Q$  plot for company B. Comparing Figs. 4(a) and 3(a), the number of incidents at company B is much higher than at company F. In addition, the variation in the number of incidents in different years is higher at company B (between ~25 and 65) than at company F (between ~0 and 15). Note that the circles on the  $Q-Q$  plot in Fig. 4(b) depart more significantly from the straight line, possibly due to the larger year-to-year variation in the number of incidents as well as the appropriateness of the of *Gamma-*

*Poisson* distribution. The circles below the straight line correspond to the safe situation where the number of incidents is less than higher than predicted, provide a warning.

The predictions in Fig. 4(b) are improved by using a *Negative Binomial* likelihood distribution with *Gamma* and *Beta* prior distributions. The prior distribution for 1990 is obtained using  $\alpha = \beta = 0.001$ , and  $a = b = 1.0$ , providing a relatively flat distribution in the region of interest; that is, a non-informative prior distribution. The *Negative Binomial* distribution provides better agreement for company B, while the *Poisson* distribution is preferred for company F.

### 3.2. Statistical analysis of incident causes and equipment types

In this analysis, for each company, Bayesian models are formulated for each cause and equipment type. Because of the large variations in the number of incidents observed

1 over the years, the performance of the *Gamma-Poisson*  
 Bayesian models differ significantly. For company F, Figs.  
 3 5(a) and (b) show the *Q-Q* plots for EF and for OE,  
 respectively. Fig. 5(a) shows better agreement with the  
 5 model because the variation in the number of incidents  
 related to EF is small, while the variation in the number of  
 7 incidents related to OE is more significant. This is  
 consistent with the expectation that equipment perfor-  
 9 mance varies less significantly than operator performance  
 over time.

11 Figs. 6(a) and (b) show the *Q-Q* plots for EF and for  
 OE, respectively, at company B. When comparing Figs.  
 13 5(a) and 6(a), the predictions of the numbers of EF at  
 company B are poorer than at company F using the  
 15 *Poisson* distribution, but are improved using the *Negative*  
*Binomial* distribution. This is similar to the predictions for  
 17 the total numbers of incidents at company B, as shown in  
 Fig. 4(b), compared with those at company F, as shown in  
 19 Fig. 3(b). Yet, the predictions for the OE are comparable at

companies F and B, and consequently, the larger variation  
 in reporting incidents at company B are attributed to the  
 larger variation in the numbers of EF.

Figs. 7(a-d) show the *Q-Q* plots for incidents associated  
 with the process units, storage vessels, heat-transfer  
 equipment, and compressors/pumps at company B using  
*Poisson* and *Negative Binomial* distributions. The *Negative*  
*Binomial* distribution is better for incidents associated with  
 the process units, compressors/pumps, and heat-transfer  
 equipment, while the *Poisson* distribution is preferred for  
 storage vessels.

3.3. Statistical analysis of chemicals involved

For each company, an attempt was made to identify  
 trends for each of the top five chemicals associated with the  
 largest number of incidents in the Harris County obtained  
 from the NRC database. However, no specific trends for a  
 particular chemical associated with a higher number of

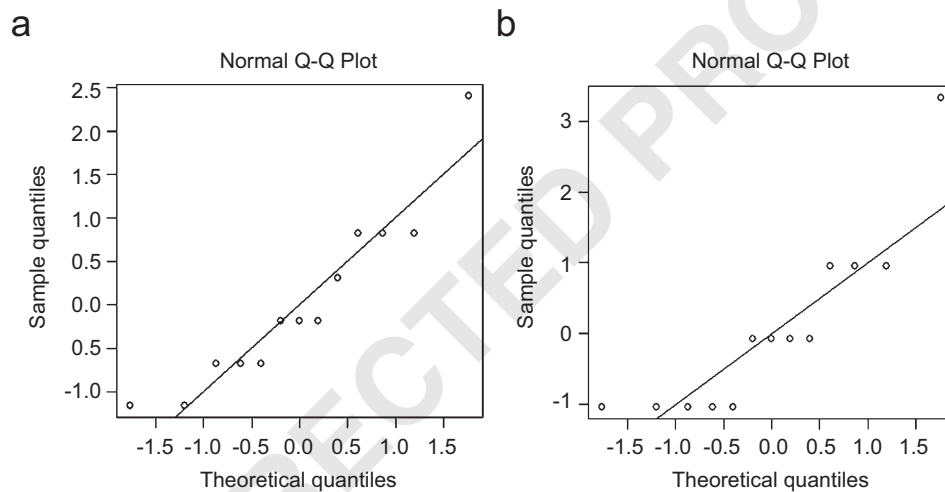


Fig. 5. Company F: (a) equipment failures, (b) operator errors.

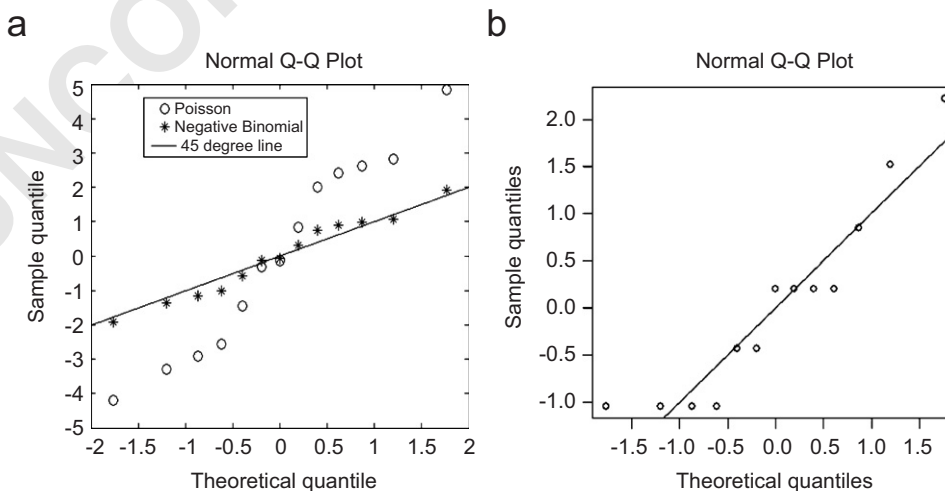


Fig. 6. Company B: (a) equipment failures, (b) operator errors.

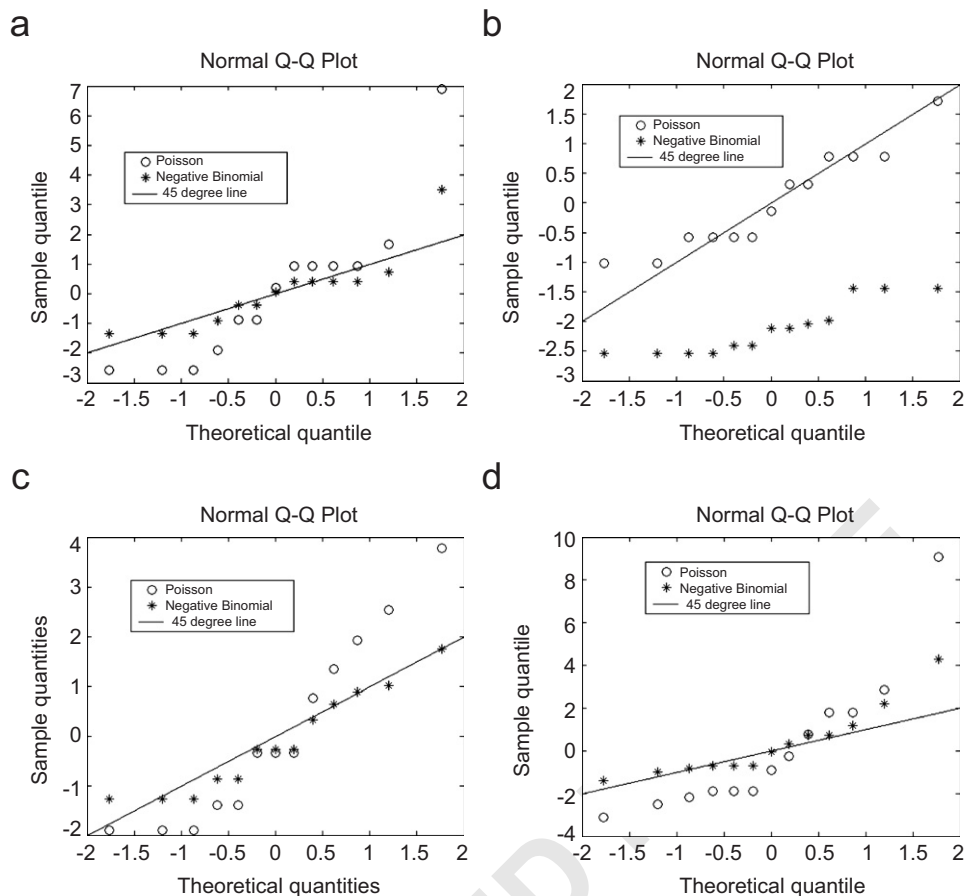


Fig. 7. Company B: (a) process units, (b) storage vessels, (c) Heat-transfer equipment, and (d) compressors/pumps.

incidents in all of the companies were observed. This could be because different products are produced in varying amounts by different companies. It might be preferable to carry out the analysis for a company that manufactures similar chemicals at different locations or for different companies that produce similar products.

### 3.4. Statistical analysis of the day of the week

For each of the seven companies, Table 2 summarizes the model checking of the Bayesian predictive distributions of the days of the week, with the mean,  $E$ , and variance,  $V$ , of  $z$  tabulated. Again, the predictions improve with the total number of incidents observed for a company. As seen, the mean and variance of  $z$  indicate that higher deviations are observed on Wednesdays and Thursdays for all of the companies, except company G. Lower deviations occur at the beginning of the week and over the weekends. To understand this observation, more information appears to be necessary; for example, (1) defining the operator shift and maintenance schedules, (2) carrying out operator surveys, (3) determining operator work loads, and (4) relating the data on the causes of the incidents to the days of the week, identifying more specific patterns. Furthermore, the higher means and variances for company G on

Friday and Saturday suggest that additional data are needed to generate a reliable Bayesian model.

### 3.5. Rates of EF and OE

In this section, for an incident, the probabilities of the involvement of each of the 13 equipment types and the probabilities of their causes (EF, OE and O) are modeled. The tree in Fig. 8 shows, for each incident, the possible causes, and for each cause, the possible equipment types. Note that alternatively the tree could show, for each incident, the possible equipment types followed by the possible causes.  $x_1, x_2, x_3$  are the probabilities of causes EF, OE, and O for an incident, and  $d_1, d_2, d_3$  are the cumulative numbers of incidents at the end of each year.  $e_1, e_2, e_3, \dots, e_{13}$  are the probabilities of the involvement of equipment types,  $E_1, E_2, \dots, E_{13}$ , in an incident through different causes, where  $M_1 + N_1 + O_1, M_2 + N_2 + O_2, M_3 + N_3 + O_3, \dots, M_{13} + N_{13} + O_{13}$  are the cumulative number of incidents associated with each equipment type.

The prior distributions of the probability of  $x_i$  are modeled using *Beta* distributions with parameters  $a_i$  and  $b_i$ :

$$f(x_i) \propto (x_i)^{a_i-1} (1-x_i)^{b_i-1}, \quad i = 1, \dots, 3, \quad (3)$$

having means  $= a_i/(a_i + b_i)$  and variances  $= a_i b_i / (a_i + b_i)^2 (a_i + b_i + 1)$ . These conjugate *Beta* prior distributions

1 Table 2  
 2 Q-Q plot properties for day of the week analysis of incidents

	Mon	Tue	Wed	Thru	Fri	Sat	Sun
5 A	0.027, 1.5	0.015, 1.06	0.032, 1.55	0.046, 1.9	0.023, 1.31	0.022, 1.23	0.055, 1.93
6 B	0.032, 1.53	0.047, 1.8	0.06, 2.12	0.058, 2.05	0.035, 1.55	0.027, 1.25	0.033, 1.46
7 C	0.027, 1.28	0.024, 1.21	0.047, 1.67	0.048, 1.62	0.031, 1.33	0.019, 1.002	0.039, 1.48
8 D	0.15, 2.3	0.165, 2.7	0.2, 2.96	0.2, 3.22	0.13, 2.44	0.126, 2.22	0.27, 3.4
9 E	0.038, 1.06	0.037, 1.19	0.086, 1.66	0.078, 1.64	0.11, 1.89	0.07, 1.46	0.036, 0.96
10 F	0.034, 1.06	0.06, 1.27	0.04, 1.08	0.87, 0.05	0.035, 0.98	0.043, 1.01	0.07, 1.22
11 G	0.06, 1.09	0.14, 1.29	0.14, 1.29	0.14, 1.29	7.84, 29.26	15.82, 58.48	0.23, 1.96

Entry in each cell-E(z), V(z)

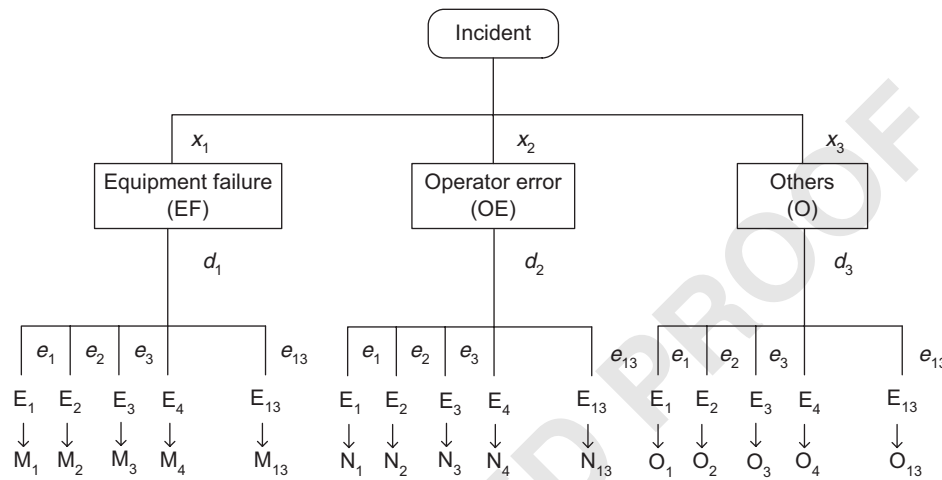


Fig. 8. Tree of causes and equipment types involved in an incident.

are updated using *Bernoulli's* likelihood distribution to obtain the posterior distribution of the probability of  $x_i$ :

$$f(x_i|Data) \propto (x_i)^{a_i-1+d_i}(1-x_i)^{b_i-1+\sum_{k=1, \neq i}^3 d_k} f(x_i). \quad (4)$$

The posterior distributions, which are also *Beta* distributions having parameters,  $a_i+d_i$ , and  $b_i+\sum_{k=1, \neq i}^3 d_k$ , change at the end of each year as  $d_i$  change.  $a_1$  and  $b_1$  are assumed to be 1.0 to give a flat, non-informative, prior distribution;  $a_2$  and  $b_2$  are assumed to be 0.998 and 1.002 to give a nearly flat, non-informative, prior distribution; and  $a_3$  and  $b_3$  are 0.001 and 0.999. Consequently, the mean prior probabilities of EF, OE, and O are 0.5, 0.499, and 0.001, respectively, as shown in Fig. 9(a).

The posterior means and variances are obtained over the years 1990–2002 for each of the seven companies. Fig. 9(a) shows the probabilities,  $x_1$ ,  $x_2$ , and  $x_3$ , of the causes EF, OE, and O for an incident at company F. Using the data at the end of each year, the probabilities increase from 0.5 for the EF, decrease from 0.499 for the OE, and increase from 0.001 for the others, with the OE approaching slightly higher values than those for the others.

Similarly, analyses for the probabilities of the equipment types,  $e_1, e_2, \dots, e_{13}$ , are carried out using *Beta* distribu-

tions,  $f(e_i)$  and  $f(e_i|Data)$ , with the *Data*,  $M_1+N_1+O_1, M_2+N_2+O_2, M_3+N_3+O_3, \dots, M_{13}+N_{13}+O_{13}$ . The prior distributions of the probabilities of  $e_i$  are modeled using *Beta* distributions with parameters  $p_i$  and  $q_i$ :

$$f(e_i) \propto (e_i)^{p_i-1}(1-e_i)^{q_i-1}, \quad i = 1, \dots, 13, \quad (5)$$

having means =  $p_i/(p_i+q_i)$  and variances =  $p_iq_i/(p_i+q_i)^2(-p_i+q_i+1)$ . These conjugate *Beta* prior distributions are updated using *Bernoulli's* likelihood distribution to obtain the posterior distributions of the probabilities of  $e_i$ :

$$f(e_i|Data) \propto (e_i)^{p_i-1+M_i+N_i+O_i} \times (1-e_i)^{q_i-1+\sum_{k=1, \neq i}^{13} M_k+N_k+O_k} f(e_i). \quad (6)$$

The posterior distributions, which are also *Beta* distributions having parameters,  $p_i+M_i+N_i+O_i$ , and  $q_i+\sum_{k=1, \neq i}^{13} M_k+N_k+O_k$ , change at the end of each year as  $M_i+N_i+O_i$  change. The parameters,  $p_i$  and  $q_i$ , are chosen to give flat, non-informative, prior distributions.

The posterior means and variances are obtained over the years 1990–2002 for each of the thirteen equipment types at each of the seven companies. Fig. 9(b) shows, for an incident, that the probability of the involvement of the PV decreases over time. Similarly, the probabilities for the



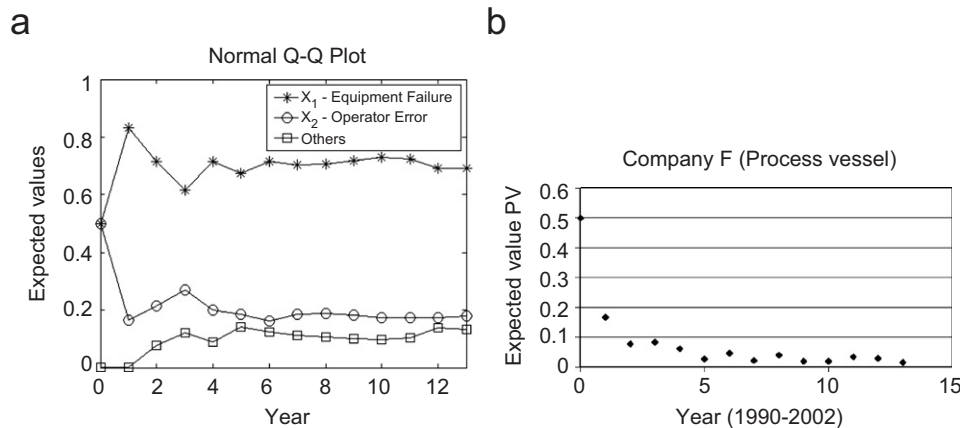


Fig. 9. Probabilities of  $x_i$  for company F: (a) EF, OE, and others, (b) PV.

Table 3  
OE/EF ratio for the petrochemical (P) and specialty chemical (S) companies

Company	A (P)	B (P)	C (S)	D (P)	E (S)	F (S)	G (S)
OE/EF ratio	0–0.3	0–0.22	0–0.75	0–0.5	0–0.667	0–0.667	0–0.5

other equipment types approach stable values after a few years with occasional departures from their mean values.

### 3.5.1. Equipment and human reliabilities

By comparing the causes of incidents between the EF and OE, insights regarding equipment and human reliabilities are obtained. In Table 3, where the range of the annual OE/EF ratio for all of the companies is shown, incidents involving EF exceed incidents involving OE. As mentioned in Section 3.1, the low OE/EF ratios are probably due to the operator bias when reporting incidents. Nevertheless, for petrochemical companies, the ratio is much lower than for specialty chemical companies. This is anticipated because the manufacture of specialty chemicals involves more batch operations, increasing the likelihood of OE.

### 3.6. Specialty chemicals and petrochemicals

To identify trends in the manufacture of specialty chemicals and petrochemicals, data for companies C, E, F, and G are combined and compared with the combined data for companies A, B, and D. Note that this is advantageous when the data for a single company are insufficient to identify trends, and when it is assumed that the lumped data for each group of companies are identically and independently distributed (i.i.d.). For these reasons, all of the analyses in Sections 3.1–3.5 were repeated with the data for specialty chemical and petrochemical manufacturers lumped together. Because the number of datum entries in each lumped data set is increased, the circles on the  $Q-Q$  plot lie closer to the

straight line. However, the cumulative predictions for the specialty chemical and petrochemical manufacturers differ significantly from those for the individual companies. Hence, it is important to carry out company specific analyses. Nevertheless, when insufficient data are available for each company, the cumulative predictions for specialty chemical and petrochemical manufacturers are preferable. Furthermore, when insufficient lumped data are available for the specialty chemicals and petrochemical manufacturers, trends may be identified by combining the data for all of the companies.

### 3.7. Modeling the loss-severity distribution using EVT

For rare events with extreme losses, it is important to identify those that exceed a high threshold. EVT is a powerful and fairly robust framework to study the tail behavior of a distribution. Embrechts et al. (1997) provide an overview of EVT as a risk management tool, discussing its potential and limitations. In another study, McNeil (1997) examines the estimation of the tails of the loss-severity distributions and the estimation of quantile risk measures for financial time-series using EVT. Herein, EVT, which uses the generalized Pareto distribution (GPD), is employed to develop a loss-severity distribution for the seven chemical companies. Other methods use the log-normal, generalized extreme value, Weibull, and Gamma distributions.

The distribution of excess values of losses,  $l$ , over a high threshold,  $u$ , is defined as:

$$F_u(y) = Pr\{l - u \leq y | l > u\} = \frac{F(y + u) - F(u)}{1 - F(u)}, \quad l \in L, \quad (7)$$

which represents the probability that the value of  $l$  exceeds the threshold,  $u$ , by less than or equal to  $y$ , given that  $l$  exceeds the threshold,  $u$ , where  $F$  is the cumulative probability distribution, and  $L$  is the set of losses. This is the so-called *loss-severity* distribution. Note that, for the NRC database,  $l$  is defined in Section 3.7.1. For sufficiently high threshold,  $u$ , the distribution function of the excess may be approximated by the GPD,  $G(l)$ , and consequently,  $F_u(y)$  converges to the GPD as the threshold becomes large. The GPD is

$$G(l) = \begin{cases} 1 - \left(1 + \xi \frac{l-u}{\beta}\right)^{-1/\xi} & \text{if } \xi \neq 0 \\ 1 - e^{-l/\beta} & \text{if } \xi = 0 \end{cases}, \quad (8)$$

where  $\beta$  is the scale parameter,  $\xi$  is the shape parameter, and the tail index is  $\xi^{-1}$ . Note that the GPD reduces to different distributions depending on  $\xi$ . The distribution of excesses may be approximated by the GPD by choosing  $\xi$  and  $\beta$  and setting a high threshold,  $u$ . The parameters of the GPD can be estimated using various techniques; for example, the maximum likelihood method and the method of probability-weighted moments.

### 3.7.1. Loss-severity distribution of the NRC database

Because few incidents have high severity levels, the incidents analyzed for the seven companies are assumed to be i.i.d. Consequently, the incidents for a specific company (internal data) are combined with those for the other companies (external data) to obtain a common *loss-severity* distribution for the seven companies. The loss associated with an incident,  $l$ , is calculated as a weighted sum of the numbers of evacuations,  $N_e$ ; injuries,  $N_i$ ; hospitalizations,  $N_h$ ; fatalities,  $N_f$ ; and damages,  $N_d$ :

$$l = w_e N_e + w_i N_i + w_h N_h + w_f N_f + w_d N_d, \quad (9)$$

where  $w_e = \$100$ ,  $w_i = \$10,000$ ,  $w_h = \$50,000$ ,  $w_f = \$2,000,000$ , and  $w_d = 1$ , with  $N_d$  reported in dollars. Note the weighting factors were adjusted to align with the company performance histories.

For the NRC database, the threshold value,  $u$ , was chosen to be \$10,000. As expected, the NRC database has few incidents that have a sizable loss. Only 157 incidents among those reported had monetary loss ( $l > 0$ ), 64 exceeded the threshold, and 108 exceeded or equaled the threshold. A software package, Extreme Value Analysis in MATLAB (EVIM) Gencay et al. (2001), obtained the parameters of the GPD,  $\xi = 0.8688$  and  $\beta = 1.7183 \times 10^4$ , for the NRC database using the maximum likelihood method. Fig. 10 shows the predictions of  $F_u(y)$ , the cumulative probability of the losses,  $l$ , that exceed or equal the threshold,  $u$ . Note that while the cumulative distribution of the losses could be improved with data from more

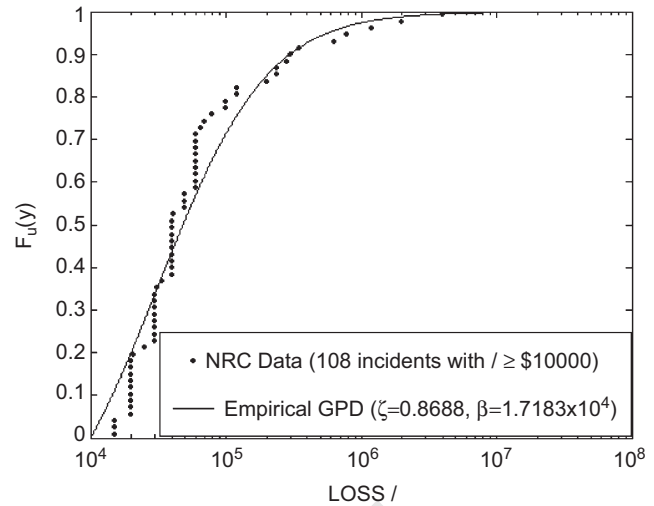


Fig. 10. Loss-severity distribution of the NRC database.

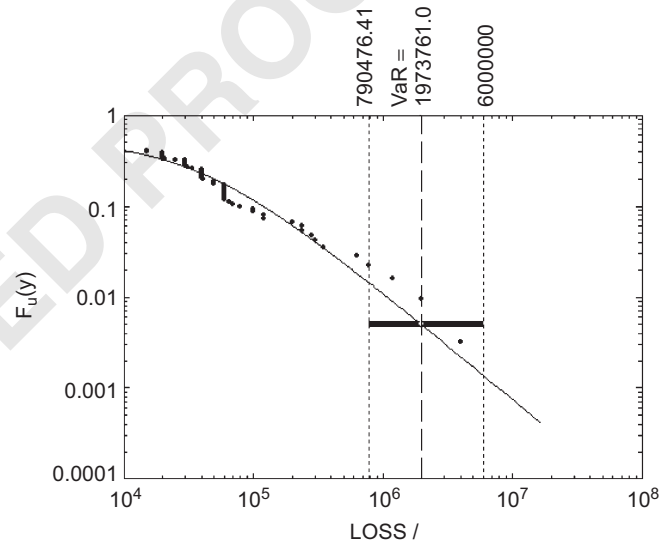


Fig. 11. Tail behavior of the *loss-severity* distribution for companies A–G.

companies in Harris County, the predictions in Fig. 10 are considered to be satisfactory.

By graphing  $\log(1 - F_u(y))$ , Fig. 11 emphasizes the tail of the *loss-severity* distribution, with the value at risk (VaR) defined at 99.5% ( $1 - F_u(y) = 0.005$ ) cumulative probability equal to  $\$1.97 \times 10^6$  and the lower and upper bounds on the 95% confidence interval equal to  $\$7.9 \times 10^5$  and  $\$6.0 \times 10^6$ , respectively. The VaR is a forecast of a specified percentile (e.g., 99.5%), usually in the right tail, of the *loss-severity* distribution over some period (e.g., annually).

## 4. Operational risk

Several types of risks, for example, credit, market, and operational risks are encountered by chemical companies. In this work, the primary focus is on calculating the

operational risk associated with a chemical company, which is defined as the risk of direct or indirect losses resulting from inadequate or failed internal resources, people, and systems, or from external events.

Capital charge (that is, CaR) of a company due to operational risk is calculated herein. Capital charge is obtained from the *total loss* distribution (to be defined below) using the VaR. Computation of the *total loss* distribution is a common statistical approach in the actuarial sciences. This paper applies this approach to risk analysis in the chemical industries. There are four methods for obtaining capital charge associated with operational risk: (i) the basic indicator approach (BIA), (ii) the standardized approach (SA), (iii) the internal measurement approach (IMA), and (iv) the loss distribution approach (LDA). The LDA (Klugman, Panjer, & Willmot, 1998) is considered to be the most sophisticated, and is used herein.

In the LDA, the annual frequency distribution of incidents is obtained using internal data, while the *loss-severity* distribution of an incident is obtained using internal and external data, as discussed in Section 3.7.1. By multiplying these two distributions, the *total loss* distribution is obtained.

Fig. 12 shows a schematic of the *total loss* distribution for a chemical company. The *expected* loss corresponds to the mean (expected) value and the *unexpected* loss is the value of the loss for a specified percentile (e.g., 99.5%) minus the *expected* loss. Note that, in some circles, the CaR is defined as the *unexpected* loss. However, herein, in agreement with other institutions, the CaR is the sum of the *expected* and *unexpected* losses, at the 99.5 percentile of the *total loss* distribution.

Highly accurate estimates of the CaR are difficult to compute due to the scarcity of internal data for the extreme events at most companies. Also, internal data are biased towards low-severity losses while external data are biased towards high-severity losses. Consequently, a mix of internal and external data is needed to enhance the statistical significance. Furthermore, it is important to balance the cost of recording very low-severity data and the truncation bias or accuracy loss resulting from the use of unduly high thresholds.

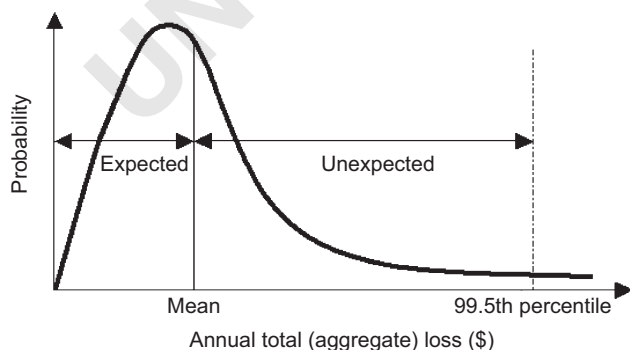


Fig. 12. Schematic of *total loss* distribution for a chemical company.

As when estimating the frequency of incidents (Section 2), a frequency distribution is obtained initially using Bayesian theory for events with losses that exceed a threshold,  $u$ . Because operational risks are difficult to estimate shortly after operations begin, conservative estimates of the parameters of the *Poisson* distribution may be obtained. In these cases, the sensitivity of the CaR to the frequency parameter should be examined. After the frequency distribution is obtained, it is multiplied with the *loss-severity* distribution and the FFT is used to calculate the *total loss* distribution.

#### 4.1. FFT algorithm

The algorithm for computing the *total loss* distribution using the FFT is described in this section. Aggregate losses are represented as the sum,  $Z$ , of a random number,  $N$ , of individual losses,  $l_1, l_2, \dots, l_N$ . The characteristic function of the total loss,  $\phi_z(t)$ , is:

$$\begin{aligned} \phi_z(t) &= E[e^{it(Z)}] = E_N[E[e^{it(l_1+l_2+\dots+l_N)}|N]] \\ &= E_N[\phi_l(t)^N] = P_N(\phi_l(t)), \end{aligned} \tag{10}$$

where  $P_N$  is the probability generating function of the frequency of incidents,  $N$ , and  $\phi_l$  is the characteristic function of the *loss-severity* distribution. The FFT produces an approximation of  $\phi_z$  and, using  $\phi_z$ , the inverse fast-Fourier transform (IFFT) gives  $f_z(Z)$ , the discrete probability distribution of the total (aggregate) loss. The details of the FFT, IFFT, and the characteristics function are found elsewhere (Klugman et al., 1998).

First,  $n_p = 2^r$  for some integer  $r$  is chosen, where  $n_p$  is the desired number of points in the distribution of total losses, such that the *total loss* distribution has negligible probability outside the range  $[0, n_p]$ . Herein,  $r = 13$  provides a sufficiently broad range. It can be adjusted according to the number of incidents in a company. The next steps in the algorithm are:

1. The *loss-severity* distribution is transformed from continuous to discrete using the method of rounding (Klugman et al., 1998). The span is assumed to be \$20,000 in line with the threshold for the GPD. The discrete loss-severity vector is represented as  $f_l = [f_l(0), f_l(1), \dots, f_l(n_p-1)]$ .
2. The FFT of the discrete loss-severity vector is carried out to obtain the characteristic function of the *loss-severity* distribution:  $\phi_l = \text{FFT}(f_l)$ .
3. The probability generating function of the frequency,  $P_N(t) = e^{\lambda(t-1)}$ , is applied, element-by-element, to the FFT of the discrete loss-severity vector to obtain the characteristic function of the *total loss* distribution:  $\phi_z = P_N(\phi_l)$ .
4. The IFFT is applied to  $\phi_z$  to recover the discrete distribution of the total losses:  $f_z = \text{IFFT}(\phi_z)$ .

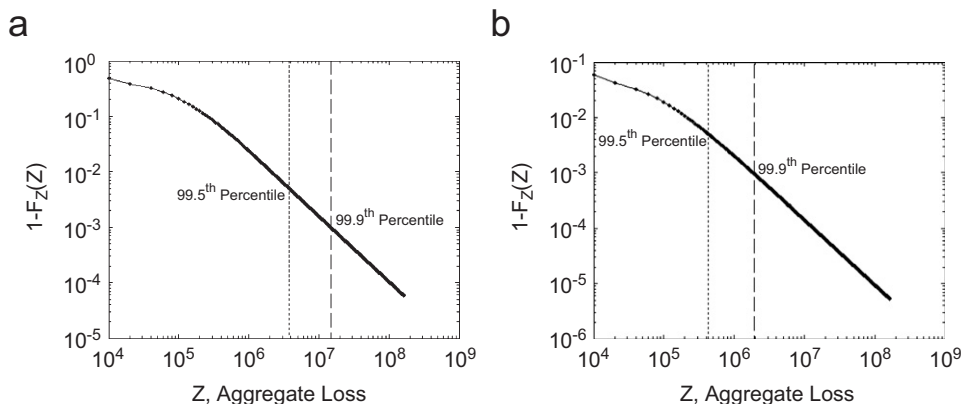


Fig. 13. Total loss distribution for: (a) company B, (b) company F.

#### 4.2. Total loss distribution for companies B and F

The *Poisson* frequency parameters for companies B and F, obtained using internal data for each company, are  $\lambda_B = 0.8461$  and  $\lambda_F = 0.0769$ . These are obtained using Bayesian theory for their incident data through the years 1 to  $n-1$  (1990–2001) for incidents having losses that exceed or equal the threshold, \$10,000. The low  $\lambda_F$  indicates the low probability of incidents having significant losses in company F. For company B,  $\lambda_B$  indicates that about one event, with  $l > \$10,000$ , is anticipated in the next year. Note that the *loss-severity* distributions in Figs. 10 and 11 are obtained using both internal and external data.

Fig. 13(a) shows the tail of the cumulative plot of the *total loss* distribution for company B. The total loss at the 99.5th percentile is  $\$3.76 \times 10^6$  and at the 99.9th percentile is  $\$14.1 \times 10^6$ . When  $\lambda_B \gg 1$ , a much higher value of CaR is expected. Similarly, Fig. 13(b) shows the tail for company F. The total loss at the 99.5th percentile is  $\$0.43 \times 10^6$  and at the 99.9th percentile is  $\$1.78 \times 10^6$ . As expected, the CaR for company F is lower than for company B by an order of magnitude.

Hence, this method provides plant-specific estimates of the CaR. Such calculations should be performed by chemical companies to provide better estimates for insurance premiums and to add quantitative support for safety audits.

#### 5. Conclusions

Statistical models to analyze accident precursors in the NRC database have been developed. They:

1. Provide Bayesian models that facilitate improved company-specific estimates, as compared with lumped estimates involving all of the specialty chemical and petrochemical manufacturers.
2. Identify Wednesday and Thursday as days of the week in which higher variations in incidents are observed.
3. Are effective for testing equipment and human reliabilities, indicating that the OE/EF ratio is lower for

petrochemical than specialty chemical companies.

4. Are beneficial for obtaining the value at risk (VaR) from the *loss-severity* distribution using EVT and the capital at risk (CaR) from the *total loss* distribution.

Consistent reporting of incidents is crucial for the reliability of this analysis. In addition, the predictive errors are reduced when: (i) sufficient incidents are available for a specific company to provide reliable means, and (ii) less variation occurs in the number of incidents from year-to-year. Furthermore, to obtain better predictions, it helps to select distributions that better represent the data, properly modeling the functionality between the mean and variance of the data.

#### Acknowledgment

The interactions and advice of Professor Paul Kleindorfer of the Wharton Risk Management and Decision Center, Wharton School, University of Pennsylvania, and Professor Sam Mannan of the Mary Kay O'Connor Process Safety Center, Texas A&M University, are appreciated. Partial support for this research from the National Science Foundation through grant CTS-0553941 is gratefully acknowledged.

#### References

- Anand, S., Keren, N., Tretter, M. J., Wang, Y., O'Connor, T. M., & Mannan, M. S. (2006). Harnessing data mining to explore incident databases. *Journal of Hazardous Material*, 130, 33–41.
- Baumont, G., Menage, F., Schneiter, J. R., Spurgin, A., & Vogel, A. (2000). Quantifying human and organizational factors in accident management using decision trees: The HORAAM method. *Reliability Engineering System Safety*, 70(2), 113–124.
- Bradlow, E. T., Hardie, B. G. S., & Fader, P. S. (2002). Bayesian inference for the negative binomial distribution via polynomial expansions. *Journal of Computational and Graphical Statistics*, 11(1), 189–201.
- CCPS (1995). Process Safety Incident Database (PSID). <<http://www.ai-che.org/CCPS/ActiveProjects/PSID/index.aspx>>.
- Chung, P. W. H., & Jefferson, M. (1998). The integration of accident databases with computer tools in the chemical industry. *Computers and Chemical Engineering*, 22, S729–S732.



- 1 Elliott, M. R., Wang, Y., Lowe, R. A., & Kleindorfer, P. R. (2004).  
 2 Environmental justice: Frequency and severity of US chemical  
 3 industry accidents and the socioeconomic status of surrounding  
 4 communities. *Journal of Epidemiology and Community Health*, 58(1),  
 5 24–30.
- 6 Embrechts, P., Kluppelberg, C., & Mikosch, T. (1997). *Modelling external  
 7 events*. Berlin: Springer.
- 8 Gencay, R., Selcuk, F., & Ulugulyagci, A. (2001). EVIM: A software  
 9 package for extreme value analysis in MATLAB. *Studies in Nonlinear  
 10 Dynamics and Econometrics*, 5(3), 213–239.
- 11 Gentleman, R., Ihaka, R., Bates, D., Chambers, J., Dalgaard, J., &  
 12 Hornik, K. (2005). The R project for Statistical Computing. <[http://  
 13 www.r-project.org/](http://www.r-project.org/)>.
- 14 Goossens, L. H. J., & Cooke, R. M. (1997). Applications of some risk  
 15 assessment techniques: Formal expert judgement and accident  
 16 sequence precursors. *Safety Science*, 26(1–2), 35–47.
- 17 Kirchsteiger, C. (1997). Impact of accident precursors on risk estimates  
 18 from accident databases. *Journal of Loss Prevention in the Process  
 19 Industries*, 10(3), 159–167.
- 20 Kleindorfer, P. R., Belke, J. C., Elliott, M. R., Lee, K., Lowe, R. A., &  
 21 Feldman, H. I. (2003). Accident epidemiology and the US chemical  
 22 industry: Accident history and worst-case data from RMP\*Info. *Risk  
 23 Analysis*, 23(5), 865–881.
- 24 Klugman, S. A., Panjer, H. H., & Willmot, G. E. (1998). *Loss Models:  
 25 From data to decisions*. Wiley series in probability and statistics. Inc.  
 26 John Wiley & Sons.
- 27 Mannan, M. S., O'Connor, T. M., & West, H. H. (1999). Accident history  
 database: An opportunity. *Environmental Progress*, 18(1), 1–6.
- McNeil, A. J. (1997). Estimating the tails of loss severity distributions  
 using extreme value theory. *ASTIN Bulletin*, 27, 117–137.
- Meel, A., & Seider, W. D. (2006). Plant-specific dynamic failure  
 assessment using Bayesian theory. *Chemical Engineering Science*, 61,  
 7036–7056.
- NRC (1990). National Response Center. <[http://www.nrc.uscg.mil/  
 29 nrchp.html](http://www.nrc.uscg.mil/nrchp.html)>.
- Phimister, J. R., Oktem, U., Kleindorfer, P. R., & Kunreuther, H. (2003).  
 30 Near-miss incident management in the chemical process industry. *Risk  
 31 Analysis*, 23(3), 445–459.
- Rasmussen, K. (1996). The experience with Major Accident Reporting  
 32 System from 1984 to 1993. *European Commission, Joint Research  
 33 Center, EUR 16341 EN*.
- RMP (2000). 40 CFR Chapter IV, Accidental Release Prevention  
 34 Requirements; Risk Management Programs Under the Clean Air  
 35 Act Section 112(r)(7); Distribution of Off-Site Consequence Analysis  
 36 Information. *Final Rule, 65 FR 48108*.
- Robert, C. P. (2001). *The Bayesian choice*. New York: Springer-Verlag.
- Sonnemans, P. J. M., & Korvers, P. M. W. (2006). Accidents in the  
 37 chemical industry: Are they foreseeable? *Journal of Loss Prevention in  
 38 the Process Industries*, 19(1), 1–12.
- Sonnemans, P. J. M., Korvers, P. M. W., Brombacher, A. C., van Beek, P.  
 39 C., & Reinders, J. E. A. (2003). Accidents, often the result of an  
 40 'uncontrolled business process'—a study in the (Dutch) chemical  
 41 industry. *Quality and Reliability Engineering International*, 19(3),  
 42 183–196.
- Spiegelhalter, D., Thomas, A., Best, N., & Lunn, D. (2003). Bayesian  
 43 inference Using Gibbs Sampling (BUGS). <[http://www.mrc-bsu.ca-  
 44 m.ac.uk/bugs/welcome.shtml](http://www.mrc-bsu.cam.ac.uk/bugs/welcome.shtml)>.
- Uth, H. J. (1999). Trends in major industrial accidents in Germany.  
 45 *Journal of Loss Prevention in the Process Industries*, 12(1), 69–73.
- Uth, H. J., & Wiese, N. (2004). Central collecting and evaluating of major  
 46 accidents and near-miss-events in the Federal Republic of Germany—  
 47 results, experiences, perspectives. *Journal of Hazardous Materials*,  
 48 111(1–3), 139–145.