# Prediction of failure probability of oil wells

**João B. Carvalho[a], Dione M. Valença[b] and Julio M. Singer[c]**

[a]*Federal University of Campina Grande*
[b]*Federal University of Rio Grande do Norte*
[c]*University of São Paulo*

**Abstract.** We consider parametric accelerated failure time models with random effects to predict the probability of possibly correlated failures occurring in oil wells. In this context, we first consider empirical Bayes predictors (EBP) based on a Weibull distribution for the failure times and on a Gaussian distribution for the random effects. We also obtain empirical best linear unbiased predictors (EBLUP) using a linear mixed model for which the form of the distribution of the random effects is not specified. We compare both approaches using data obtained from an oil-drilling company and suggest how the results may be employed in designing a preventive maintenance program.

## 1 Introduction

The productivity of oil wells depends on the performance of a sub-surface equipment system that may fail with the presence of sand, corrosion, internal pressure variation, etc. One or more failures may occur during a well's lifetime and in these cases, losses may be large, since expensive corrective maintenance procedures must be considered while production is interrupted. Preventive maintenance programs (substitution of certain parts, cleaning, lubrication, etc.) may be implemented to reduce such losses. These programs, however, depend on the selection of wells with higher failure probabilities which may be estimated via statistical models.

As an example, we consider a study of the time between failures of sub-surface equipment of a sample of oil wells obtained from an oil-drilling company, between January 2000 and December 2006. Failure is defined as the complete stop in the operation of the well caused by any problem in the sub-surface equipment. Since each well may have several failures, we expect the time intervals between them to be correlated. Using a reliability terminology and considering each well as a complex repairable system, we refer to recurrent events (failures) in repairable systems. Furthermore, it is necessary to take into account the censoring that arises either because some wells are disabled from production or because of the termination of the study. For interesting literature reviews on recurrent events in reliability, we refer

to Ascher and Feingold (1984), Lawless and Thiagarajah (1996), Lugtigheid et al. (2004) and Percy and Alkali (2007).

A convenient approach to handle this type of correlated survival data (see, e.g., Hougaard (2000)) is via models with random effects as in Robinson (1991). Some authors recommend using accelerated failure time (AFT) models (see Hougaard et al. (1994) and Keiding et al. (1997)) within each unit, assuming that the omission of some important covariates may be the cause of the within-unit correlations and that the inclusion of (nonobservable) random effects may take this omission into account. In such models, the logarithm of the event times follows a linear regression on the covariate vector, so that the random effects act multiplicatively on the event times. This approach is considered in Lambert et al. (2004) or in Bolfarine and Valença (2005), for example. In this setup, a natural parametric approach, successfully applied in a variety of disciplines is to use Weibull regression models, perhaps the most widely used parametric model in survival analysis and reliability experiments (see Lawless (2003)).

Lambert et al. (2004) employ empirical Bayes methods (see, e.g., Carlin and Louis (1998)) to analyze kidney transplant data under such models and use the mode of the posterior distribution of the realized random effects as the predictor. In their analysis, different combinations of the distribution of the random effects and lifetimes are considered. In this paper, assuming a Weibull distribution for the time intervals between failures, we consider a linear mixed model approach, where only the existence of the first two moments of the distribution of the random effects is required and use empirical best linear unbiased predictors (EBLUP) as an alternative. We compare the conditional probabilities of failure obtained via empirical Bayes and EBLUP approaches in an analysis of the oil well data described above.

The model and the two alternative prediction approaches are described in Section 2. Data analysis is described in Section 3. Results are compared in Section 4 and a brief discussion is presented in Section 5.

## 2 Accelerated failure time model with random effects

Let $n_i$ observations be recorded on the $i$th of $k$ units along the duration of the study. Let $T_{ij}$ denote the time between the $(j-1)$th and the $j$th failure of the $i$th unit, $i = 1, \ldots, k$, $j = 1, \ldots, n_i$, not including repair time. The accelerated failure time model with random effects is

$$\ln T_{ij} = b_i + \boldsymbol{\beta}^\top \mathbf{x}_{ij} + \sigma \varepsilon_{ij}, \tag{2.1}$$

where $\mathbf{x}_{ij}$ denotes a $p \times 1$ vector of covariates with the first component equal to 1, $\boldsymbol{\beta}$ represents a $p \times 1$ vector of fixed (but unknown) parameters, $\sigma$ is a scale parameter, $b_i$, $i = 1, \ldots, k$ are independent and identically distributed unobserved random variables (random effects) with null means and common variance $\sigma_b^2$ and $\varepsilon_{ij}$, $i = 1, \ldots, k$, $j = 1, \ldots, n_i$ are independent and identically distributed unobserved

random errors with common and known mean and variance $\sigma_\varepsilon^2$. Furthermore, we assume that $\mathrm{Cov}(b_i, \varepsilon_{ij}) = 0$, $i = 1, \ldots, k$, $j = 1, \ldots, n_i$. This model reduces to the usual accelerated failure time model when $\sigma_b^2 = 0$ (see, e.g., Lawless (2003)).

Our objective is to estimate the conditional probability that a selected unit fails in an interval of length $\Delta t$ given that it has functioned correctly for at least $t$ units of time, for example,

$$P(t < T_{ij} \leq t + \Delta t | T_{ij} > t, b_i, \mathbf{x}_{ij}) = \frac{S(t|b_i, \mathbf{x}_{ij}) - S(t + \Delta t | b_i, \mathbf{x}_{ij})}{S(t|b_i, \mathbf{x}_{ij})}, \quad (2.2)$$

where $S(t|b_i, \mathbf{x}_{ij}) = P(T_{ij} > t | b_i, \mathbf{x}_{ij})$ is the conditional survival function of $T_{ij}$ given $b_i$. We assume a Weibull distribution for the failure time and therefore

$$S(t|b_i, \mathbf{x}_{ij}) = \exp\{-\exp[-\sigma^{-1}(b_i + \mathbf{x}_{ij}^\top \boldsymbol{\beta})] t^{\sigma^{-1}}\}. \quad (2.3)$$

Note that in this case the variance of random errors of model (2.1) is known ($\sigma_\varepsilon^2 = \pi^2/6$).

Empirical Bayes methods can be employed to estimate (2.3); the reader is referred to Aalen and Husebye (1991) or Lambert et al. (2004) for details. To describe the procedure, we first note that because of censoring, we do not observe $\ln T_{ij}$ in all cases; instead, we observe $Y_{ij} = \min(\ln T_{ij}, \ln C_{ij})$, $i = 1, \ldots, k$, $j = 1, \ldots, n_i$, where $C_{ij}$ denotes the censoring time for the $j$th observation of the $i$th unit. Therefore, the responses are represented by $(Y_{ij}, \delta_{ij})$, where $\delta_{ij} = I(T_{ij} \leq C_{ij})$ is an indicator of failure.

The method considers the estimation of the unknown parameters of model (2.1), namely, $\boldsymbol{\lambda} = (\boldsymbol{\beta}^\top, \sigma, \sigma_b^2)^\top$, via the maximization of the marginal distribution of the responses $(Y_{ij}, \delta_{ij})$, $i = 1, \ldots, k$, $j = 1, \ldots, n_i$,

$$L(\boldsymbol{\lambda}) = \prod_{i=1}^{k} \int L_i(\boldsymbol{\beta}, \sigma | b_i) g(b_i; \sigma_b^2) \, db_i, \quad (2.4)$$

where $g(b_i; \sigma_b^2)$ is the prior density function of $b_i$ and

$$L_i(\boldsymbol{\beta}, \sigma | b_i) = \prod_{j=1}^{n_i} f(y_{ij}|b_i, \mathbf{x}_{ij})^{\delta_{ij}} S(y_{ij}|b_i, \mathbf{x}_{ij})^{(1-\delta_{ij})} \quad (2.5)$$

is the conditional likelihood function for the response of the $i$th unit given $b_i$, with $f$ and $S$, respectively denoting the conditional density function and the conditional survival function of $\ln T_{ij}$ given $b_i$. Given the Weibull assumption for $T_{ij}$, this has an extreme value distribution, where for $j = 1, \ldots, n_i$, $i = 1, \ldots, k$

$$f(y_{ij}|b_i, \mathbf{x}_{ij}) = \frac{1}{\sigma} \exp\left[\frac{y_{ij} - (b_i + \boldsymbol{\beta}^\top \mathbf{x}_{ij})}{\sigma} - \exp\left(\frac{y_{ij} - (b_i + \boldsymbol{\beta}^\top \mathbf{x}_{ij})}{\sigma}\right)\right]$$

and

$$S(y_{ij}|b_i, \mathbf{x}_{ij}) = \exp\left[-\exp\left(\frac{y_{ij} - (b_i + \boldsymbol{\beta}^\top \mathbf{x}_{ij})}{\sigma}\right)\right].$$

Letting $\widehat{\boldsymbol{\lambda}}$ denote the maximum likelihood estimator of $\boldsymbol{\lambda}$, the Bayes empirical predictor of $\mathbf{b} = (b_1, \ldots, b_k)^\top$ is the mode of the posterior distribution

$$\pi(\mathbf{b}|\mathbf{y}, \widehat{\boldsymbol{\lambda}}) = \frac{\prod_{i=1}^k L_i(\widehat{\boldsymbol{\beta}}, \widehat{\sigma}|b_i) g(b_i; \widehat{\sigma}_b^2)}{L(\widehat{\boldsymbol{\lambda}})}. \tag{2.6}$$

Given that $\pi(\mathbf{b}|\mathbf{y}, \widehat{\boldsymbol{\lambda}})$ depends on $\mathbf{b}$ only through $L_i(\widehat{\boldsymbol{\beta}}, \widehat{\sigma}|b_i)$ and $g(b_i; \widehat{\sigma}_b^2)$, the Bayes empirical predictor of $b_i$ corresponds to the point that maximizes

$$L_i(\widehat{\boldsymbol{\beta}}, \widehat{\sigma}|b_i) g(b_i; \widehat{\sigma}_b^2).$$

Depending on the form of the distribution assumed for the random effects, the integral in (2.4) may not be analytically tractable. This is the case when the assumption of normality, commonly employed for the random effects, is adopted. In such cases, several approaches have been used (see Aitkin (1999) or Breslow and Clayton (1993), e.g.). Here, an adapted Gaussian quadrature algorithm is used to approximate the integral and the maximization of the likelihood is based on an iterative method. For details, the reader is referred to Liu and Pierce (1994).

Alternatively, we may consider a linear mixed model for the analysis of the data. However, because of censoring, each component $Y_{ij}$ of the response vector $\mathbf{Y}$ has the form

$$Y_{ij} = \delta_{ij} \ln T_{ij} + (1 - \delta_{ij}) \ln C_{ij}, \tag{2.7}$$

inducing an underestimation of the true logarithm of time between failures. To bypass this problem under a nonparametric approach, we propose an adaptation of the method considered in Ageel (2002) to impute censored observations in Weibull-regression models. The idea is to replace $C_{ij}$ in (2.7) with an estimate $\widehat{C}_{ij}$ of $E(T_{ij}|T_{ij} > C_{ij}, \mathbf{x}_{ij})$, and take

$$Y_{ij}^* = \delta_{ij} \ln T_{ij} + (1 - \delta_{ij}) \ln \widehat{C}_{ij}$$

as the response variable.

Letting $\mathbf{Y}^* = (\mathbf{Y}_1^{*\top}, \ldots, \mathbf{Y}_k^{*\top})^\top$, with $\mathbf{Y}_i^* = (Y_{i1}^*, \ldots, Y_{in_i}^*)^\top$, $i = 1, \ldots, k$, the linear mixed model is

$$\mathbf{Y}^* = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \mathbf{e}, \tag{2.8}$$

where $\mathbf{X} = (\mathbf{X}_1^\top, \ldots, \mathbf{X}_k^\top)^\top$, $\mathbf{X}_i = (x_{i1}, x_{i2}, \ldots, x_{in_i})^\top$, $\mathbf{Z} = \bigoplus_{i=1}^k \mathbf{1}_{n_i}$, $\mathbf{b} = (b_1, \ldots, b_k)^\top$, $E(\mathbf{b}) = \mathbf{0}$, $\text{Var}(\mathbf{b}) = \sigma_b^2 \mathbf{I}_k$, with $\mathbf{I}_k$ denoting the identity matrix of dimension $k$ and $\bigoplus_{i=1}^k \mathbf{a}_i$ representing the direct sum of the vectors $\mathbf{a}_i$. Furthermore, $\mathbf{e}$ is a $n \times 1$ vector of random errors, with $n = \sum_{i=1}^k n_i$, $\mathbf{e} = \sigma[\varepsilon - E(\varepsilon)]$, $\text{Var}(\mathbf{e}) = (\sigma^2 \pi^2/6)\mathbf{I}_n$ and uncorrelated with $\mathbf{b}$. These definitions imply

$$E(\mathbf{Y}^*) = \mathbf{X}\boldsymbol{\beta}, \tag{2.9}$$

$$\text{Var}(\mathbf{Y}^*) = \mathbf{V} = \sigma_b^2 \mathbf{Z}\mathbf{Z}^\top + (\sigma^2 \pi^2/6)\mathbf{I}_n \tag{2.10}$$

$$\text{Cov}(\mathbf{b}, \mathbf{Y}^{*\top}) = \mathbf{C} = \sigma_b^2 \mathbf{Z}^\top. \tag{2.11}$$

When $\mathbf{C}$ and $\mathbf{V}$ are known, the best linear unbiased estimator (BLUE) of $\boldsymbol{\beta}$ and the best linear unbiased predictor (BLUP) of $\mathbf{b}$ are obtained as the solutions to the well known Henderson equations (Henderson (1975)) and are respectively given by

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V}^{-1} \mathbf{Y}^* \quad \text{and} \quad \tilde{\mathbf{b}} = \mathbf{C} \mathbf{V}^{-1} (\mathbf{Y}^* - \mathbf{X} \widehat{\boldsymbol{\beta}}). \qquad (2.12)$$

In practice, $\sigma_b^2$ and $\sigma^2$ are unknown and must be replaced by estimates in (2.12), generating the empirical best linear unbiased estimator (EBLUE) and the empirical best linear unbiased predictor (EBLUP). The most common methods of estimation of the variance components are maximum likelihood and restricted maximum likelihood, but this requires the specification of the form of the distribution of $\mathbf{e}$ and $\mathbf{b}$. For details, the reader is referred to Jiang (1997) among others. Non-parametric estimators based on quadratic functions of the data, like the minimum norm quadratic unbiased estimator (MINQUE) or the minimum variance quadratic unbiased estimator (MIVQUE) considered in Rao (1970, 1971a, 1971b) may be employed when the form of the underlying distribution is not specified. Other methods of estimation are described in Searle et al. (1992) or Demidenko (2004). The analysis of mixed linear models with censored observations in a parametric setup is considered in Hughes (1999) and Pettitt (1986), among others.

In our context, prediction of random effects under accelerated failure time models is particularly appealing; it is computationally simpler than the empirical Bayes approach and does not require an assumption on the form of the distribution of the random effects.

## 3 Data analysis

To identify the oil wells with the highest probabilities of failure, a total of 2374 failure times, of which 563 (23.7%) were censored, was recorded for 616 oil wells. For the purpose of this study, five covariates were included, namely,

- Production (PROD) in $m^3$/day;
- Elevation method: mechanical pumping (MP) or progressive cavity (PC);
- Age at failure (AGE) in years;
- Region: RA, RB, RC and RD;
- Depth of the oil pump (DEPTH) in m.

Note that regions and elevation methods are represented by dummy variables in the model. An initial exploratory analysis (not shown) followed by the variable selection strategy recommended by Collett (1994) based on likelihood ratio tests suggested the model

$$\ln T_{ij} = b_i + \beta_0 + \beta_{\text{prod}} PROD_{ij} + \beta_{\text{bm}} MP + \beta_{\text{age}} AGE_{ij} + \beta_{\text{rb}} RB + \beta_{\text{rc}} RC$$
$$+ \beta_{\text{rd}} RD + \beta_{\text{depth}} DEPTH_i + \beta_{\text{prod}*\text{rb}} PROD_{ij} * RB$$

$$+ \beta_{\mathrm{prod}*\mathrm{rc}} PROD_{ij} * RC + \beta_{\mathrm{prod}*\mathrm{rd}} PROD_{ij} * RD$$

$$+ \beta_{\mathrm{depth}*\mathrm{rb}} DEPTH_i * RB + \beta_{\mathrm{depth}*\mathrm{rc}} DEPTH_i * RC$$

$$+ \beta_{\mathrm{depth}*\mathrm{rd}} DEPTH_i * RD + \sigma \varepsilon_{ij},$$

where $b_i \sim N(0, \sigma_b^2)$, and $T_{ij}$ follows a Weibull distribution conditionally on $b_i$, $i = 1, \ldots, 616$, $j = 1, \ldots, n_i$.
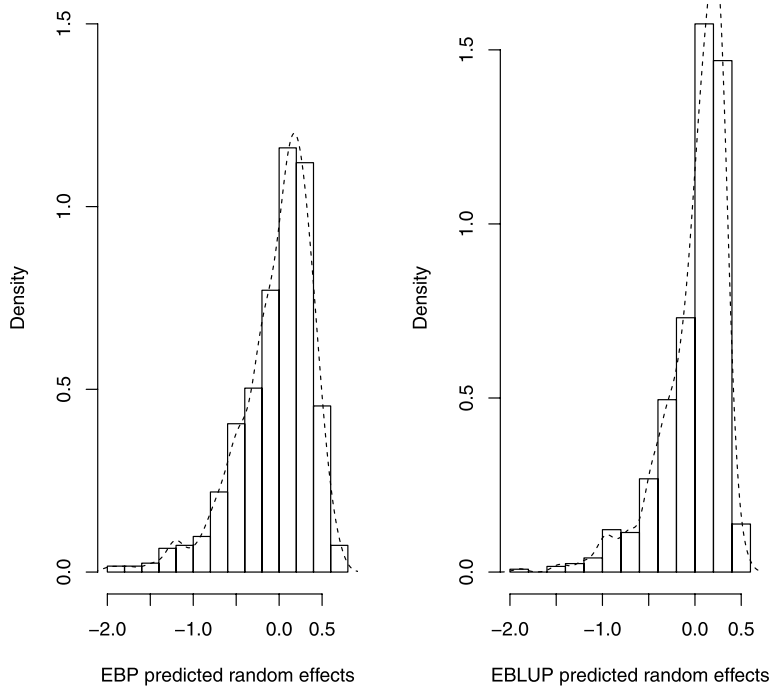
The model was fitted via empirical Bayes methods using PROC NLMIXED in SAS (SAS Institute (2009), version 9.1). Initial values for $\boldsymbol{\beta}$ and $\sigma$ were obtained from fitting a standard accelerated failure time model, for example, in which no random effects are included and the initial value for $\sigma_b^2$ was set to 0.1.

The EBLUE of the parameters and the EBLUP of the random effects were obtained from (2.12) with the variance components estimated via MINQUE methods. Computations were conducted in R (R Development Core Team (2010), version 2.12.0). The EBLUEs with corresponding standard errors are displayed in Table 1, both for the complete data (2000–2006) and for the data corresponding to the period 2000–2005. This last option was considered with the purpose of model validation.

Histograms for the predicted random effects along with kernel density estimates (see, e.g., Scott (1992)) are presented in Figure 1. They suggest that the distribution of the random effects is asymmetric, so that the Gaussian assumption may not be appropriate. In this sense, the EBLUP might be a better option since their deriva-

**Table 1** *Estimates and standard errors of the parameters of model*

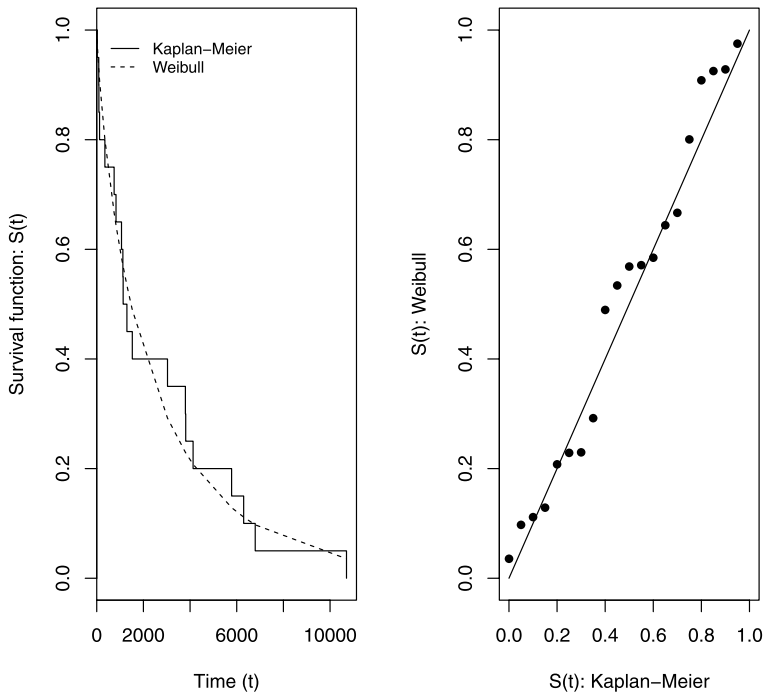| Parameter | 2000 to 2006 | | | | 2000 to 2005 | | | |
|---|---|---|---|---|---|---|---|---|
| | Empirical Bayes | | EBLUE | | Empirical Bayes | | EBLUE | |
| | Estimate | S.E. | Estimate | S.E. | Estimate | S.E. | Estimate | S.E. |
| $\beta_0$ | 7.4698 | 0.2572 | 7.1274 | 0.2480 | 7.2389 | 0.2071 | 6.5580 | 0.2679 |
| $\beta_{\mathrm{prod}}$ | −0.0498 | 0.0091 | −0.0298 | 0.0115 | −0.0228 | 0.0087 | −0.0056 | 0.0121 |
| $\beta_{\mathrm{bm}}$ | 0.5271 | 0.1284 | 0.2804 | 0.1278 | 0.4970 | 0.1057 | 0.4048 | 0.1384 |
| $\beta_{\mathrm{age}}$ | 0.0787 | 0.0077 | 0.0514 | 0.0071 | 0.0465 | 0.0064 | 0.0425 | 0.0082 |
| $\beta_{\mathrm{rb}}$ | 1.3428 | 0.3071 | 0.8513 | 0.2994 | 1.3649 | 0.2825 | 1.4458 | 0.3486 |
| $\beta_{\mathrm{rc}}$ | 0.9589 | 0.2275 | 0.8952 | 0.2182 | 0.9511 | 0.1744 | 1.2200 | 0.2260 |
| $\beta_{\mathrm{rd}}$ | 1.8259 | 0.3279 | 1.3441 | 0.3188 | 1.1991 | 0.2685 | 1.0125 | 0.3454 |
| $\beta_{\mathrm{depth}}$ | 0.0021 | 0.0004 | 0.0014 | 0.0004 | 0.0019 | 0.0003 | 0.0018 | 0.0004 |
| $\beta_{\mathrm{prod}*\mathrm{rb}}$ | 0.0323 | 0.0153 | 0.0347 | 0.0184 | 0.0171 | 0.0153 | 0.0008 | 0.0200 |
| $\beta_{\mathrm{prod}*\mathrm{rc}}$ | 0.0189 | 0.0139 | 0.0002 | 0.0174 | 0.0073 | 0.0132 | −0.0102 | 0.0181 |
| $\beta_{\mathrm{prod}*\mathrm{rd}}$ | −0.0407 | 0.0188 | −0.0467 | 0.0225 | −0.0308 | 0.0177 | −0.0384 | 0.0237 |
| $\beta_{\mathrm{rb}*\mathrm{depth}}$ | −0.0020 | 0.0004 | −0.0013 | 0.0005 | −0.0017 | 0.0005 | −0.0016 | 0.0006 |
| $\beta_{\mathrm{rc}*\mathrm{depth}}$ | −0.0028 | 0.0005 | −0.0022 | 0.0005 | −0.0025 | 0.0004 | −0.0029 | 0.0006 |
| $\beta_{\mathrm{rd}*\mathrm{depth}}$ | −0.0028 | 0.0006 | −0.0020 | 0.0006 | −0.0022 | 0.0005 | −0.0019 | 0.0007 |
| $\sigma$ | 1.2001 | 0.0253 | 1.1592 | – | 1.1320 | 0.0258 | 1.1789 | – |
| $\sigma_b^2$ | 0.4369 | 0.0582 | 0.3587 | – | 0.1292 | 0.0311 | 0.2160 | – |

**Figure 1**   *Histograms for the Bayes empirical predicted random effects and EBLUP.*

tion only requires the existence of the first and second moments of the underlying random variables.

To evaluate the appropriateness of the conditional Weibull distribution assumption, we fitted Weibull models for 20 wells that experienced twelve or more failures during the study period. The fitted model as well as the Kaplan Meyer survival curve for one of these wells along with the corresponding QQ plot are presented in Figure 2. Similar plots were constructed for all the selected wells and are not shown because they exhibit the same type of behaviour. Although the analyses were based on a small number of observations, they show no strong evidence against the Weibull assumption.

## 4  Comparison of empirical Bayes and EBLUP

Based on the ideas of Harrell Jr. and Frank (2001), we fitted the model using the 2000–2005 data (see Table 1) and used it to predict the conditional failure probabilities in $\Delta t$ hours. These quantities were computed according to (2.2)–(2.3) using either the empirical Bayes approach or the linear mixed model approach to obtain the predicted random effects. The value considered for $t$ in the expression for the conditional failure probabilities, was the actual operation time (in hours) between
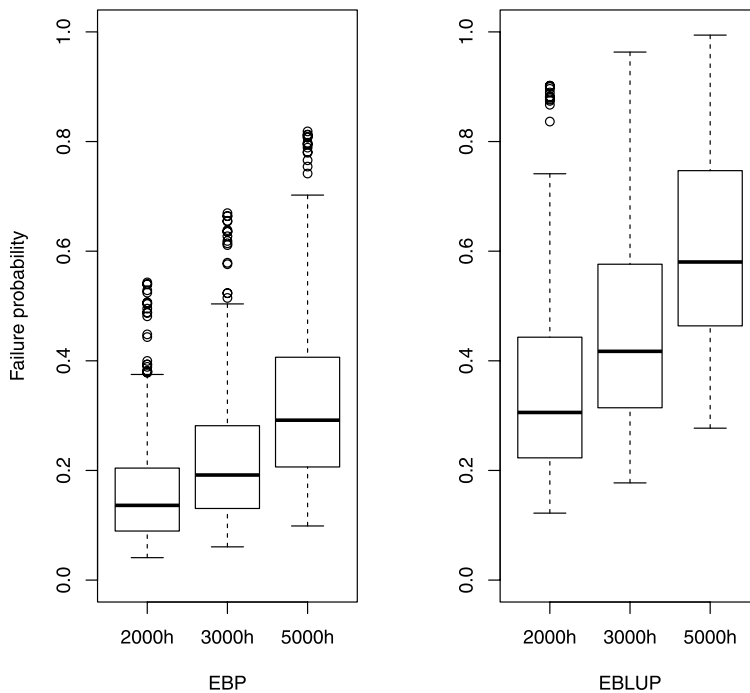
**Figure 2** *Diagnostics to evaluate the appropriateness of the Weibull model.*

the last repair in 2005, and 12/31/2005. We evaluated the ability of the model in indicating right decisions with respect to preventive maintenance in a given $\Delta t$ time interval. We considered the following steps to validate the model.

(i) We predicted the conditional failure probabilities for $\Delta t = 2000$ h, 3000 h and 5000 h, for 105 wells that failed in 2006. The corresponding box-plots are displayed in Figure 3. These plots suggest that the failure probabilities increase nonlinearly with time and may help to develop a preventive maintenance policy based on the identification of the wells with the largest predicted conditional failure probabilities in $\Delta t$ hours.

(ii) We assumed that a well should undergo preventive maintenance when the conditional failure probability predicted was greater than an arbitrary but fixed cut-off point $p_0$. The optimum cut-off point may be obtained from a ROC curve (see, e.g., Zweig and Campbell (1993)). ROC curves based on failure probabilities in $\Delta t = 3000$ h (median time of failure) predicted via EBP and EBLUP are presented in Figure 4, with several cut-off points indicated in the main diagonal. The optimum cut-off point is the one located closest to the point with coordinates (0, 1). With the available data, it is difficult to decide for an optimal cut-off point since the corresponding ROC curve is close to the no-discrimination line.
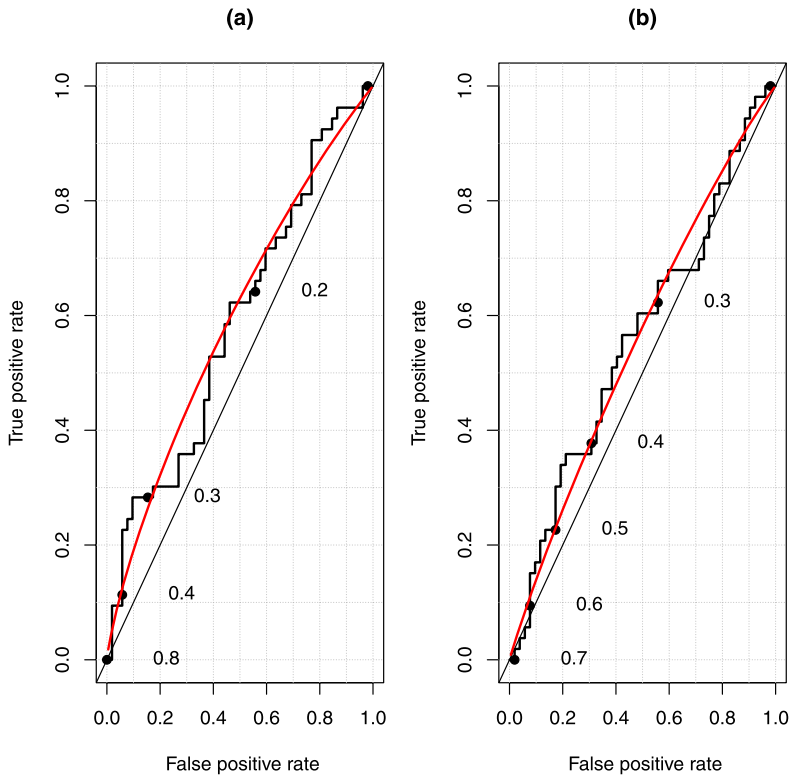
**Figure 3** *Box plots for the conditional failure probabilities (EBP and EBLUP).*

(iii) We predicted the failure probability within $\Delta t = 3000$ h from January 1, 2006 for each of the selected wells using cut-off points $p_0 = 0.4$ and $p_0 = 0.3$ and compared the results with the observed failure status in 2006. We evaluated whether the decision based on (ii) was correct (if preventive maintenance was indicated and a failure occurred or preventive maintenance was not indicated and in fact no failure was observed) or not (if preventive maintenance was indicated and no failure occurred or preventive maintenance was not indicated and in fact a failure was observed).

(iv) In each case, we computed the sensitivity, specificity, positive predictive value (PPV), negative predictive value (NPV) and accuracy of the procedure (see, e.g., Hosmer and Lemeshow (2000)). The results are summarized in Tables 2 and 3.

Accuracy, NPV and PPV of the decision rule are similar under both procedures, while EBP shows better specificity and EBLUP, better sensitivity.

## 5 Discussion

We considered two approaches to predict conditional failure probabilities of oil wells based on possibly correlated data. Both consider an accelerated failure time

**Figure 4**  *ROC curves for the conditional failure probabilities via EBP* (a) *and EBLUP* (b).

**Table 2**  *Decisions and true failure status for wells in* 2006

| | | Preventive maintenance decision | | | |
| | | Empirical Bayes | | EBLUP | |
| $P_0$ | Failure in 3000 h | yes | no | yes | no |
|---|---|---|---|---|---|
| 0.3 | yes | 15 | 38 | 34 | 19 |
| | no | 8 | 44 | 29 | 23 |
| 0.4 | yes | 7 | 46 | 20 | 33 |
| | no | 3 | 49 | 16 | 36 |

model with random effects and assume a Weibull distribution for the time between failures. While the empirical Bayes approach is usually based on a Gausssian distribution for the random effects, the linear mixed model approach is based only on the existence of the second moment of that distribution and is computationally simpler. A full hierarchical Bayesian model could be considered instead; although it

**Table 3**  *Accuracy of decision approaches*

| $p_0$ | Accuracy measure | Empirical Bayes | EBLUP |
|---|---|---|---|
| 0.3 | Sensitivity | 0.28 | 0.64 |
|  | Specificity | 0.85 | 0.44 |
|  | PPV | 0.65 | 0.54 |
|  | NPV | 0.54 | 0.55 |
|  | Accuracy | 0.56 | 0.54 |
| 0.4 | Sensitivity | 0.13 | 0.38 |
|  | Specificity | 0.94 | 0.69 |
|  | PPV | 0.70 | 0.56 |
|  | NPV | 0.52 | 0.52 |
|  | Accuracy | 0.53 | 0.53 |

may be derived under different assumptions for the random effects distribution, it is computationally more intensive than the empirical Bayesian approach we adopted. Furthermore, to take advantage of the Bayesian paradigm, it would require the elicitation of the prior distribution of $\lambda$ which would require information that, for confidential reasons, we could not access.

The linear mixed model accommodates censored observations by imputing the mean time of failure. However, according to simulations in Ageel (2002), the procedure has no significant gain with respect to the analysis based on the original data if censoring affects less than 30% of the observations. In fact, in Section 3, we have not obtained significant changes in predictions via EBLUP when we adopted this procedure in data analysis in which 23% of the observations were censored in the complete data and only 5% were censored in the period 2000 to 2005.

Our intent here was to outline statistical methodology that may be employed to obtain predictors of failure probabilities. Unfortunately we did not have access to other covariates that could have been useful to generate more accurate results.

The predicted conditional failure probabilities may be employed in the implementation of a preventive maintenance policy. This, however depends on the cut-off point $p_0$ that in turn, depends on the relative costs of corrective and preventive measures. Given such costs, we may optimize the choice of this cut-off point.

Future research in this area may be directed to (i) relaxing the Weibull assumption and adopting a Cox regression approach, (ii) improvement of the imputation procedure required in the linear mixed model approach, (iii) developing residual analysis for the linear mixed model in this setup, (iv) considering a full hierarchical Bayes approach to obtain the predictors and (v) considering linear mixed models under elliptically symmetric or skew normal distributions.

## Acknowledgments

## References

Aalen, O. and Husebye, E. (1991). Statistical analysis of repeated events forming renewal processes. *Statistics in Medicine* **10**, 1227–1240.

Ageel, M. I. (2002). A novel means of estimating quantiles for 2-parameter Weibull distribution under the right random censoring model. *Journal of Computational and Applied Mathematics* **149**, 373–380. MR1937288

Aitkin, M. (1999). A general maximum likelihood analysis of variance components in generalized linear models. *Biometrics* **55**, 117–128. MR1705676

Ascher, H. and Feingold, H. (1984). *Repairable Systems Reliability: Modeling, Inference, Misconceptions and Their Causes*. New York: Marcel Dekker. MR0762088

Bolfarine, H. and Valença, D. (2005). Testing homogeneity in Weibull-regression models. *Biometrical Journal* **47**, 707–720. MR2209066

Breslow, N. E. and Clayton, D. G. (1993). Approximate inference in generalized linear mixed models. *Journal of the American Statistical Association* **88**, 9–25.

Carlin, B. P. and Louis, T. A. (1998). *Bayes and Empirical Bayes Methods for Data Analysis*. New York: Chapman & Hall. MR1427749

Collett, D. (1994). *Modelling Survival Data in Medical Research*. New York: Chapman & Hall.

Demidenko, E. (2004). *Mixed Models: Theory and Applications*. New York: Wiley. MR2077875

Harrell Jr., F. E. and Frank E. (2001). *Regression Modeling Strategies: With Applications to Linear Models, Logistic Regression and Survival Analysis*. New York: Springer.

Henderson, C. R. (1975). Best linear unbiased estimation and prediction under a selection model. *Biometrics* **31**, 387–396.

Hosmer, D. W. and Lemeshow, S. L. (2000). *Applied Logistic Regression*, 2nd ed. New York: Wiley.

Hougaard, P. (2000). *Analysis of Multivariate Survival Data*. New York: Springer. MR1777022

Hougaard, P., Myglegaard, P. and Borch-Johnsen, K. (1994). Heterogeneity models of disease susceptibility, with application to diabetic nephropathy. *Biometrics* **50**, 1178–1188.

Hughes, J. P. (1999). Mixed effects models with censored data with application to HIV RNA levels. *Biometrics* **55**, 625–629.

Jiang, J. (1997). Asymptotic properties of the empirical BLUP and BLUE in mixed linear models. *Statistica Sinica* **8**, 861–885. MR1651513

Keiding, N., Andersen, P. K. and Klein, J. P. (1997). The role of frailty models and accelerated failure time models in describing heterogeneity due to omitted covariates. *Statistics in Medicine* **16**, 215–224.

Lambert, P., Collett, D., Kimber, A. and Johnson, R. (2004). Parametric accelerated failure time models with random effects and an application to kidney transplant survival. *Statistics in Medicine* **23**, 3177–3192.

Lawless, J. F. (2003). *Statistical models and methods for lifetime data*, 2nd ed. New York: Wiley. MR1940115

Lawless, J. F. and Thiagarajah, K. (1996). A point-process model incorporating renewals and time trends, with application to repairable systems. *Technometrics* **38**, 131–138.

Liu, Q. and Pierce, D. A. (1994). A note on Gauss–Hermite quadrature. *Biometrika* **81**, 624–629. MR1311107

Lugtigheid, D., Banjevic, D. and Jardine, A. K. S. (2004). Modelling repairable system reliability with explanatory variables and repair and maintenance actions. *IMA Journal of Management Mathematics* **15**, 89–110. MR2069534

Percy, D. and Alkali, B. (2007). Scheduling preventive maintenance of oil pumps using generalised proportional intensities models. *International Transactions in Operational Research* **14**, 547–563.

Pettitt, N. A. (1986). Censored observations, repeated measures and mixed effects models: An approach using EM algorithm and normal errors. *Biometrika* **73**, 635–643. MR0897855

R Development Core Team (2010). *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. Available at http://www.R-project.org.

Rao, C. R. (1970). Estimation of heteroscedastic variances in linear models. *Journal of the American Statistical Association* **65**, 161–172. MR0286221

Rao, C. R. (1971a). Estimation of variance and covariance components—MINQUE theory. *Journal of Multivariate Analysis* **1**, 257–275. MR0301869

Rao, C. R. (1971b). Minimum variance quadratic unbiased estimation of variance components. *Journal of Multivariate Analysis* **1**, 445–456. MR0301870

Robinson, G. K. (1991). That BLUP is a good thing: The estimation of random effects. *Statistical Science* **6**, 15–51. MR1108815

SAS Institute (2009). *SAS/STAT User's Guide, Version 9.1*. Cary: SAS Institute.

Scott, D. W. (1992). *Multivariate Density Estimation. Theory, Practice and Visualization*. New York: Wiley. MR1191168

Searle, S. R., Casella, G. and McCulloch, C. E. (1992). *Variance Components*, 2nd ed. New York: Wiley. MR1190470

Zweig, M. H. and Campbell, G. (1993). Receiver-operating characteristic (ROC) plots: A fundamental evaluation tool in clinical medicine. *Clinical Chemistry* **39**, 561–577.

J. B. Carvalho
Unidade Acadêmica de Matemática e Estatística
Universidade Federal de Campina Grande
R. Aprígio Veloso, Bairro Universitário
Campina Grande, RN, 58429-900
Brasil
E-mail: joaobc@dme.ufcg.edu.br

D. M. Valença
Departamento de Estatística
Universidade Federal do Rio Grande do Norte
Av. Senador Salgado Filho, Lagoa Nova
Natal, RN, 59078-970
Brasil
E-mail: dione@ccet.ufrn.br

J. M. Singer
Departamento de Estatística
Universidade de São Paulo
Rua do Matão, 1010
São Paulo, SP, 05508-090
Brasil
E-mail: jmsinger@ime.usp.br