

Default Estimation and Expert Information: All Likely Dataset Analysis and Robust Validation

Nicholas M. Kiefer

Departments of Economics and Statistical Sciences,

Cornell University,

490 Uris Hall, Ithaca, NY 14853-7601, US.

email: nmk1@cornell.edu

June, 2006

Abstract

Default is a rare event, even in segments in the midrange of a bank's portfolio. Most loans are in the middle-risk categories and sample sizes are often high. Expert information, crucial in inference about defaults in low-default segments, continues to be important in moderate-default situations. The method of All Likely Datasets, based on sufficient statistics and expert information, is used to characterize likely datasets for analysis. Robustness is illustrated with an ϵ -mixture of priors.

1 Introduction

Estimation of default probabilities (PD), loss given default (LGD, a fraction) and exposure at default (EAD) for portfolio segments containing reasonably homogeneous assets is essential to prudent risk management. It is also, not coincidentally, crucial for compliance with Basel II rules for banks using the IRB approach to determine capital requirements (Basel Committee on Banking Supervision 2004). Estimation of small probabilities is tricky, and I will focus on estimating PD. The focus here is on segments in the middle of the risk profile of the portfolio. Although the risk is in the middle of the asset mix, the probability of default is still "small."

We will see that it is likely to be about 0.01; defaults, though seen, are rare. The bulk of a typical bank's commercial loans are concentrated in these segments (they may differ across banks). Very low risk institutions are relatively few in number and they have access to capital through many avenues in addition to commercial loans. Very high risk investments are largely avoided and when present are often due to the reclassification of a safer loan as conditions change. To put this in perspective, the middle-quality loans are approximately S&P Baa or Moody's BBB. Of course, the bulk of these loans are to unrated companies and the bank has done their own rating to assign the loans to risk "buckets." The focus of this paper is on estimation of the default probability for such a risk bucket on the basis of historical information and expert knowledge. We do not analyze a specific data set. These are typically proprietary and results are necessarily specific to that dataset. Instead, we introduce the "All Likely Data" (ALD) approach. We use the theory of sufficient statistics to define dataset types characterized by the number of defaults for a particular sample size. The number of types is linear in the sample size, while the number of distinct datasets is exponential. This affords considerable simplification. Next, we use expert information to identify likely types, and then run analyses for all likely types - a set of types corresponding to the most likely datasets. Since defaults are expected to be rare events, a small number of types characterize the likely samples. Finally, we conduct a robustness analysis, in the spirit of validation exercises required of banks under Basel II.

2 The Characterization of Uncertainty

There are a number of arguments that uncertainty is best described in terms of probabilities. In the case of default modeling, where measuring and controlling risk is the aim, it is natural to focus on anticipating defaults, or at least anticipating the aggregate number of defaults. Suppose there are a number of default configurations, labeled E_i , E_1 might be that asset 1 and only asset 1 defaults. E_k might be a more complicated event, like "assets 3, 4 and 22-30 default." You wish to assign numbers to these events and to use these numbers to describe the likelihood of the events. Let us adopt the convenient notation that $E_i = 1$ if event E_i occurs and $E_i = 0$ if not. Since you are interested in prediction, let us consider choosing numbers x_i to

minimize your forecast error:

$$s(x_1, \dots, x_n | E_1, \dots, E_n) = \sum_{i=1}^n (x_i - E_i)^2$$

and this should be minimized for whatever configuration of E occurs. If there are n assets, there are 2^n configurations of events, each with its own sum of squares. Clearly, there are many possible solutions. It is possible to characterize the feasible solutions and describe some properties implied by this simple structure. Since the events are unknown and we wish to describe the uncertainty about these events we will call the descriptions predictions. Some authors object to this terminology since prediction may (reasonably) be interpreted to apply to something in the future, while we use it to apply to anything unknown, future or not. The suggested alternative term, prevision, suffers the same potential confusion. Consequently we use prediction, noting the extended application to current, past or unrealized unknown events.

Definition 1 *Coherence:* A set of numbers, predictions, $\{x_i^*\}$ is coherent if there is no alternative set of predictions $\{x_i\}$ such that

$$\begin{aligned} s(x_1, \dots, x_n | E_1, \dots, E_n) &\leq s(x_1^*, \dots, x_n^* | E_1, \dots, E_n) \quad \forall \{E_i\} \text{ configurations} \\ s(x_1, \dots, x_n | E_1, \dots, E_n) &< s(x_1^*, \dots, x_n^* | E_1, \dots, E_n) \quad \text{for some } \{E_i\} \end{aligned}$$

Coherence is a basic requirement for a set of predictions to be acceptable. It requires that there not be an alternative set of predictions which do at least equally well in all configurations of events and better in at least one. An economist will recognize the condition as Pareto optimality, a decision theorist as admissibility. This requirement is enough to insure that the predictions must combine like probabilities. The first implication is convexity:

Theorem 1 *Convexity:* $0 \leq x \leq 1$.

Proof. Suppose $x < 0$. This x only appears as $(x - 1)^2$ and x^2 . Both of these can be reduced by increasing x to 0. The same logic establishes that $x \leq 1$. ■

Next we have additivity:

Theorem 2 *Additivity:* Let x refer to the event E and y the event $\sim E$. Then $x + y = 1$.

Proof. Again, it suffices to consider only terms involving x and y , as the other terms entering the scores will be the same in all cases. The isoscore sets corresponding to the event E are spheres centered on $(1,0)$ in the x, y plane, and the sets corresponding to $\sim E$ are spheres centered on $(0,1)$. The coherent choices occur at tangencies, which lie on the line segment connecting the centers of the spheres (at any other point, both scores can be reduced; for example by moving toward this line segment along a path perpendicular to the segment). Thus, $x + y = 1$. Simple conditioning in the definition of events can be used to obtain the corollary that if x is the event E , y the event F , and z the event E or F , and E and F are mutually exclusive, then $x + y = z$. ■

Finally, we can establish the multiplication rule:

Theorem 3 *Multiplication:* Let x correspond to E , y to F given E , and z to E and F . Then $z = xy$.

Proof. There are 3 configurations to consider: EF , in which case the score is $(x - 1)^2 + (y - 1)^2 + (z - 1)^2$; $E(1 - F)$, giving $(x - 1)^2 + y^2 + z^2$, and $(1 - E)(1 - F)$, giving $x^2 + z^2$. The isoscore sets are thus the spheres centered on $(1,1,1)$ and $(1,0,0)$ and the cylinders centered on the y axis in the (x, y, z) coordinate system. Consequently, the coherent triplets (x, y, z) must lie in the plane containing $(1,1,1)$, $(0,y,0)$ and $(1,0,0)$. Thus we can write (using also convexity), $(x, y, z) = \alpha(1, 1, 1) + \beta(0, y, 0) + (1 - \alpha - \beta)(1, 0, 0)$ and solving we see that this requires $z = xy$. ■

These three properties are often taken as defining a system of probabilities. This development using coherence is due to (De Finetti 1974). In fact, the approach is more general; the quadratic specification is inessential to give probability as the coherent measure of uncertainty (Lindley 1982b).

Of course, the probability approach to describing and modeling uncertainty is central to risk management and to the requirements of Basel II. There is no serious argument that the probability approach is wrong or inappropriate for modeling uncertain future defaults as well as other unknowns. The fact that probabilities combine in accordance with convexity, additivity and multiplication is central for moving from probabilities of default on an asset, to default rates in a segment, to rates in a portfolio, and to a default probability for the bank. Economists do not need convincing that probabilistic reasoning is appropriate for modeling. The coherence argument is sketched here since it is less well accepted that uncertainty

about the unknown default probability can be usefully modeled in exactly the same way as uncertainty about unknown defaults, for exactly the same reasons.

3 A Statistical Model for Defaults

The simplest and most common probability model for defaults of assets in a homogeneous segment of a portfolio is the Binomial, in which the defaults are independent across assets and over time, and defaults occur with common probability θ . Note that specification of this model requires expert judgement, that is, information. Denote the expert information by e . The role of expert judgement is not usually explicitly indicated at this stage, so it is worthwhile to point to its contribution. First, consider the statistical model. The independent Bernoulli model is not the only possibility. Certainly independence is a strong assumption and would have to be considered carefully. Second, are the observations really identically distributed? Perhaps the default probabilities differ across assets. Can this be modeled, perhaps on the basis of asset characteristics? The usual requirements demand an annual default probability, estimated over a sample long enough to cover a full cycle of economic conditions. Thus the probability should be marginal with respect to external conditions. For specificity we will continue with the Binomial specification. Let d_i indicate whether the i th observation was a default ($d_i = 1$) or not ($d_i = 0$). The Bernoulli model (a single Binomial trial) for the distribution of d_i is $p(d_i|\theta, e) = \theta^{d_i}(1 - \theta)^{1-d_i}$. Let $D = \{d_i, i = 1, \dots, n\}$ denote the whole data set and $r = r(D) = \sum_i d_i$ the count of defaults. Then the joint distribution of the data is

$$\begin{aligned} p(D|\theta, e) &= \prod \theta^{d_i}(1 - \theta)^{1-d_i} \\ &= \theta^r(1 - \theta)^{n-r} \end{aligned} \tag{3.1}$$

As a function of θ for given data D this is the likelihood function $L(\theta|D, e)$. Since this distribution depends on the data D only through r (n is regarded as fixed), the sufficiency principle implies that we can concentrate attention on the distribution of r

$$p(r|\theta, e) = \binom{n}{r} \theta^r (1 - \theta)^{n-r} \tag{3.2}$$

a Binomial(n, θ) distribution. This is so well known that it is easy to underappreciate the simplification obtained by passing from 3.1 to 3.2. Instead of separate

treatment for each of the 2^n datasets possible, it is sufficient to restrict attention to $n+1$ data set types, characterized by the value of r . This theory of types can be made the basis of a theory of asymptotic inference (Cover and Thomas 1991). In our application, the set of likely values of r is small, and an analysis can be done for each of these values of r , rather than for the $\binom{n}{r}$ distinct datasets corresponding to each value of r . Thus, by analyzing a few likely data set types, we analyze in effect all of the most likely data realizations. We refer to this approach as the method of all likely datasets, or ALD.

Regarded as a function of θ for fixed r , 3.2 is the likelihood function. Figure 1 shows the likelihood functions for $n=500$, our reference data set size, and $r=\{0,2,4,6,8\}$.

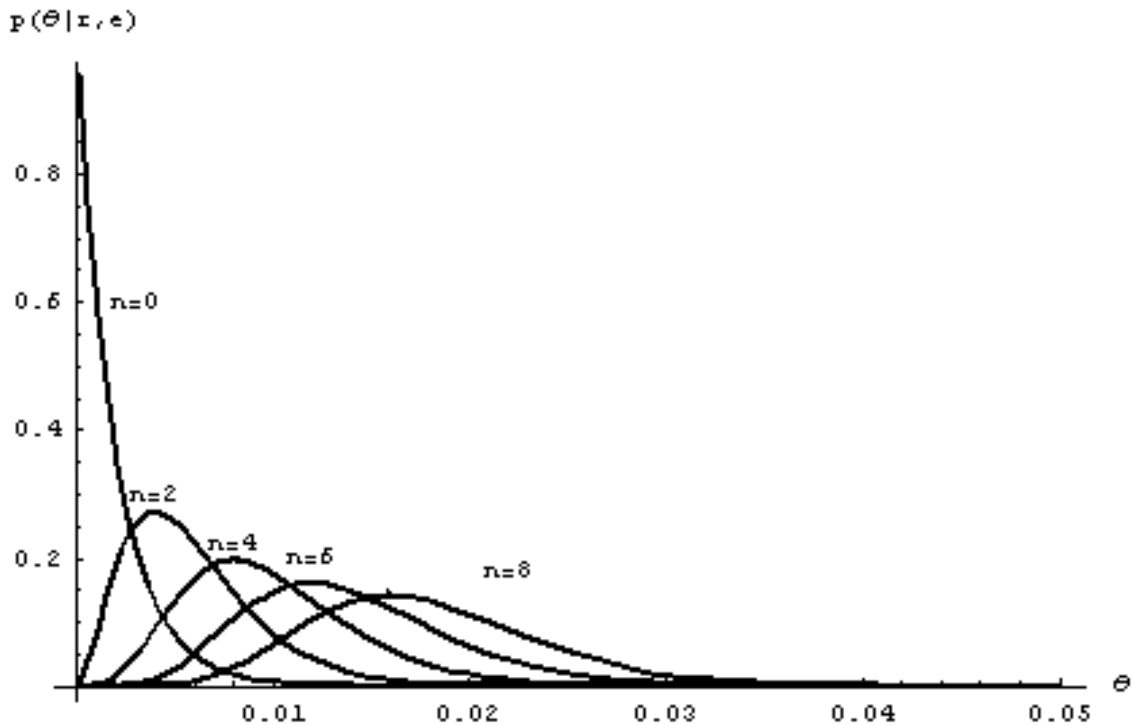


Figure 1: Likelihood Functions

4 Uncertain Default Probabilities

3.2 is a statistical model. It generates probabilities for all default configurations as a function of a single parameter θ , which remains unspecified. The default

probability θ is an unknown, but that doesn't mean that nothing is known about its value. In fact, defaults are widely studied and risk managers, modelers, validators, and supervisors have detailed knowledge on values of θ for particular portfolio segments. The point is that θ is unknown in the same sense that the future default status of a particular asset is unknown. The fact that default is in the future is not important; the key is that it is unknown and the uncertainty can be described and quantified. We have seen how uncertain defaults can be modeled. The same methods can be used to model the uncertainty about θ . Define events E_i relevant to describing the uncertainty about θ , for example $E_1 = "\theta < .0001"$; $E_2 : "\theta < .0005,"$ etc. Applying the theorem on scoring, we see that uncertainty about values of θ are coherently described by probabilities. We assemble these probability assessments into a distribution describing the uncertainty about θ given the expert information e , $p(\theta|e)$.

Now, $p(\theta|e)$ can be a quite general specification, reflecting in general the assessments of uncertainty in an infinity of possible events. This is in contrast with the case of default configurations, in which there are only a finite (though usually large) number of possible default configurations. However, this should not present an insurmountable problem. Note that we are quite willing to model the large number of probabilities associated with the possible different default configurations with a simple statistical model; in fact, a 1-parameter model. This involves an independence assumption, among other assumptions, but it simplifies the analysis and allows progress along empirical lines. The same can be done with the prior specification. That is, we can fit a few probability assessments by an expert to a suitable functional form and use that distribution to model prior uncertainty. Of course, there is some approximation involved, and care is necessary. In this regard, the situation is no different from that present in likelihood specification.

A convenient functional form is the beta distribution

$$p(\theta|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \quad (4.1)$$

which has mean $\alpha/(\alpha + \beta)$ and variance $(a - b)^2 \alpha \beta / ((\alpha + \beta)^2 (1 + \alpha + \beta))$. The special case of $\alpha = \beta = 1$ is the uniform distribution on the unit interval. This is unlikely to represent information about default probabilities, since it assigns equal probabilities to each equal length interval in $[0,1]$, but it is of great historical interest and is in common use as representing complete absence of information (it has maximal

entropy among distributions on $[0,1]$). It will be useful in constructing a robust prior for a validation step in the analysis.

A particularly easy generalization is to specify the support $\theta \in [a, b] \subset [0, 1]$. It is possible that some applications would require the support of θ to consist of the union of disjoint subsets of $[0, 1]$, but this seems fanciful in the current application. A simple starting point is the uniform $p(\theta|e) = 1/(b - a)$. This prior would again sometimes be regarded as "uninformative," since it assigns equal probability to equal length subsets of $[a, b]$. Of course, it is informative in that it requires $\theta \in [a, b]$. The mean of this distribution is $(b - a)/2$. We may think that this specification is too restrictive, in that consideration might require that intervals near the most likely value should be more probable than intervals near the endpoints. A somewhat richer specification is the beta distribution 4.1 modified to have support $[a, b]$. Let t have the beta distribution and change variables to $\theta(t) = a + (b - a)t$ with inverse function $t(\theta) = (\theta - a)/(b - a)$ and jacobian $dt(\theta)/d\theta = 1/(b - a)$. Then

$$p(\theta|\alpha, \beta, a, b) = \frac{\Gamma(\alpha + \beta)}{(b - a)\Gamma(\alpha)\Gamma(\beta)} ((a - \theta)/(a - b))^{\alpha-1} ((\theta - b)/(a - b))^{\beta-1} \quad (4.2)$$

over the range $\theta \in [a, b]$. This distribution has mean $E\theta = (b\alpha + a\beta)/(\alpha + \beta)$, allowing substantially more flexibility than the uniform. Examples of this distribution on the range $[0, 0.3]$ are graphed in Figure 2.

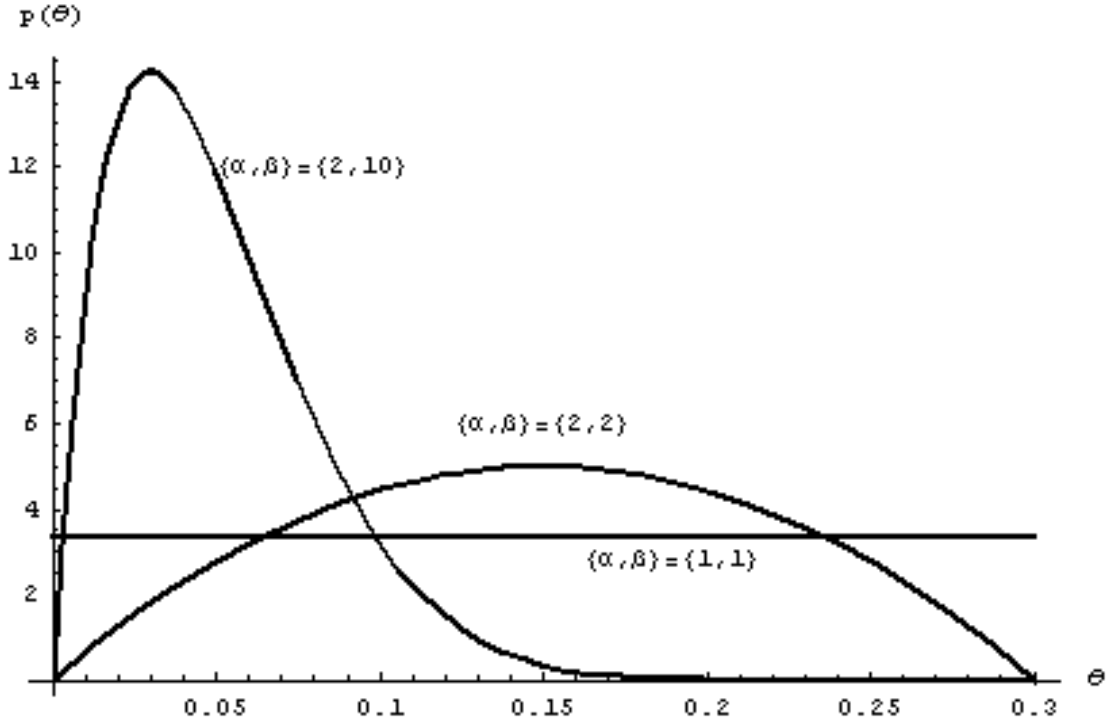


Figure 2: Beta Distributions on $[0,0.3]$

The four parameter Beta distribution allows flexibility within the range $[a,b]$, but in some situations it may be too restrictive. For example, it is unimodal. This is unlikely to be a problem for representing the prior uncertainty of an individual expert, but it may not be flexible enough to allow combination of information from many experts. A simple generalization is the 7-parameter mixture of two 4-parameter Betas with common support. The additional parameters are the two new $\{\alpha, \beta\}$ parameters and the mixing parameter λ .

$$\begin{aligned}
 p(\theta|\alpha_1, \beta_1, \alpha_2, \beta_2, a, b) &= \frac{\lambda\Gamma(\alpha_1 + \beta_1)}{(b-a)\Gamma(\alpha_1)\Gamma(\beta_1)} \left(\frac{a-\theta}{a-b}\right)^{\alpha_1-1} \left(\frac{\theta-b}{a-b}\right)^{\beta_1-1} \\
 &+ \frac{(1-\lambda)\Gamma(\alpha_2 + \beta_2)}{(b-a)\Gamma(\alpha_2)\Gamma(\beta_2)} \left(\frac{a-\theta}{a-b}\right)^{\alpha_2-1} \left(\frac{\theta-b}{a-b}\right)^{\beta_2-1}
 \end{aligned}$$

Computations with this mixture distribution are not substantially more complicated than computations with the 4-parameter Beta alone. If necessary, more mixture components with new parameters can be added, although it seems unlikely that expert information would be detailed and specific enough to require this complicated a representation. A useful further generalization is given by the 9-parameter

mixture allowing different supports for the two mixture components. The prior family is then

$$\begin{aligned}
p(\theta|\alpha_1, \beta_1, \alpha_2, \beta_2, a, b, c, d) = & \\
& \frac{I(\theta \in [a, b])\lambda\Gamma(\alpha_1 + \beta_1)}{(b - a)\Gamma(\alpha_1)\Gamma(\beta_1)}((a - \theta)/(a - b))^{\alpha_1-1}((\theta - b)/(a - b))^{\beta_1-1} \\
& + \frac{I(\theta \in [c, d])(1 - \lambda)\Gamma(\alpha_2 + \beta_2)}{(d - c)\Gamma(\alpha_2)\Gamma(\beta_2)}((c - \theta)/(c - d))^{\alpha_2-1}((\theta - d)/(c - d))^{\beta_2-1} \quad (4.3)
\end{aligned}$$

Here $[c, d]$ is the support set for the second mixture component and $I[x] = 1$ if condition x is true, 0 if false. As above, more than two mixture components could be added as needed, possibly with different support sets. By choosing enough Beta-mixture terms the approximation of an arbitrary continuous prior $p(\theta|e)$ for a Bernoulli parameter can be made arbitrarily accurate (Diaconis and Ylvisaker 1985)

5 Inference

Given the distribution $p(\theta|e)$, we can multiply the probabilities in accord with the multiplication rule to obtain the joint distribution of r , the number of defaults, and θ :

$$p(r, \theta|e) = p(r|\theta, e)p(\theta|e)$$

from which we obtain using the addition rule the marginal (predictive) distribution of r ,

$$p(r|e) = \int p(r, \theta|e)d\theta \quad (5.1)$$

If the value of the parameter θ is of main interest (rather than the prediction of the number of defaults) we can divide to obtain the conditional (posterior) distribution of θ :

$$p(\theta|r, e) = p(r|\theta, e)p(\theta|e)/p(r|e) \quad (5.2)$$

which is Bayes rule. Since Basel II places more emphasis on the default probability than on the number of defaults in a given portfolio segment, we focus our discussion on $p(\theta|r, e)$.

6 Prior Distribution

I have asked an expert to specify a portfolio and give me some aspects of his beliefs about the unknown default probability. The portfolio consists of loans that might be in the middle of a bank's portfolio. These are typically commercial loans, mostly to unrated companies. If rated, these might be about S&P Baa or Moody's BBB. The method included a specification of the problem and some specific questions followed by a discussion. General discussions of the elicitation of prior distributions are given by (Kadane, Dickey, Winkler, Smith, and Peters 1980), (Garthwaite, Kadane, and O'Hagan 2005) and (Kadane and Wolfson 1998). An example assessing a prior for a Bernoulli parameter is (Chaloner and Duncan 1983). Chaloner and Duncan follow Kadane et al in suggesting that assessments be done not directly on the probabilities concerning the parameters, but on the predictive distribution. That is, questions should be asked about observables, to bring the expert's thoughts closer to familiar ground. In the case of a Bernoulli parameter and a 2-parameter beta prior, Chaloner and Duncan suggest first eliciting the mode of the predictive distribution for a given n (an integer), then assessing the relative probability of the adjacent values. Graphical feedback is provided for refinement of the specification. Examples consider $n=20$. The suggestion to interrogate experts on what they would expect to see in data, rather than what they would expect of parameter values, is appealing and I have to some extent pursued this with our expert. This approach may be less attractive in the case of large sample sizes and small probabilities, and in our particular application, where the experts are sophisticated about probabilities. Our expert found it easier to think in terms of the probabilities directly than in terms of defaults in a hypothetical sample.

The sample period should be currently relevant, but should include a cycle, so that it is marginal with respect to business conditions. It could be argued that a recent period including the 2001-2002 period of mild downturn covers a modern cycle. A period that included the 1980's would yield higher default probabilities but these are probably not currently relevant. The default probability of interest is the current and immediate future value, not a guess at what past estimates might

be. The precise definition of default is also at issue. In the economic theory of the firm, default occurs when debt payments are missed and ownership and control of the firm passes from existing owners (shareholders in the case of a corporation) to debtholders. As a lesser criterion, loans that are assigned to "nonaccrual" may be considered defaulted. We simply note the importance of using consistent definitions in the assessment of expert information and in data definition.

We did the elicitation assuming a sample of 500 asset/years. For our application, we also considered a "small" sample of 100 observations and a "large" sample of 1000 observations, and occasionally an enormous sample of 10000 observations. Considering first the predictive distribution on 500 observations, the modal value was five defaults. Upon being asked to consider the relative probabilities of five or four defaults, conditional on four or five defaults occurring (the conditioning does not matter here, for the probability ratio, but it is thought to be easier to think about when posed in this fashion), the expert expressed some trepidation as it is difficult to think about such rare events. Ultimately, the expert gave probability ratios not achievable by the binomial model even with known probability. The expert was quite happy in thinking about probabilities over probabilities however. This may not be so uncommon in this technical area, as practioners are accustomed to working with probabilities. The mean value was 0.01. The minimum value for the default probability was 0.0001 (one basis point). The expert reported that a value above 0.035 would occur with probability less than 10%, and an absolute upper bound was 0.3. The upper bound was discussed: the expert thought probabilities in the upper tail of his distribution were extremely unlikely, but he did not want to rule out the possibility that the rates were much higher than anticipated (prudence?). Quartiles were assessed by asking the expert to consider the value at which larger or smaller values would be equiprobable given the value was less than the median, then given the value was more than the median. The median value was 0.01. The former was 0.0075. The latter, the .75 quartile, was assessed at .0125. The expert seemed to be thinking in terms of a normal distribution, perhaps using informally a central limit theorem combined with long experience with this category of assets.

This set of answers is more than enough information to determine a 4-parameter Beta distribution. I used a method of moments to fit parametric probability statements to the expert assessments. The moments I used were squared differences relative to the target values, for example $((a - 0.0001)/0.0001)^2$. The support points were quite well-determined for a range of $\{\alpha, \beta\}$ pairs at the assessed values

$\{a, b\} = [0.0001, 0.3]$. These were allowed to vary but the optimization routine did not change them beyond the 7th decimal place. Further, changing the weights did not matter much either. Probably this is due to the fact that there is almost no probability in the upper tail, so changing the upper bound made almost no difference in the assessed probabilities. Thus the rather high (?) value of b reflects the long tail apparently desired by the expert. The $\{\alpha, \beta\}$ parameters were rather less well-determined (the sum of squares function was fairly flat) and I settled on the values (‘7.9, 224.8) as best describing the expert’s information. The resulting prior distribution $p(\theta|e)$ is graphed in Figure 3.

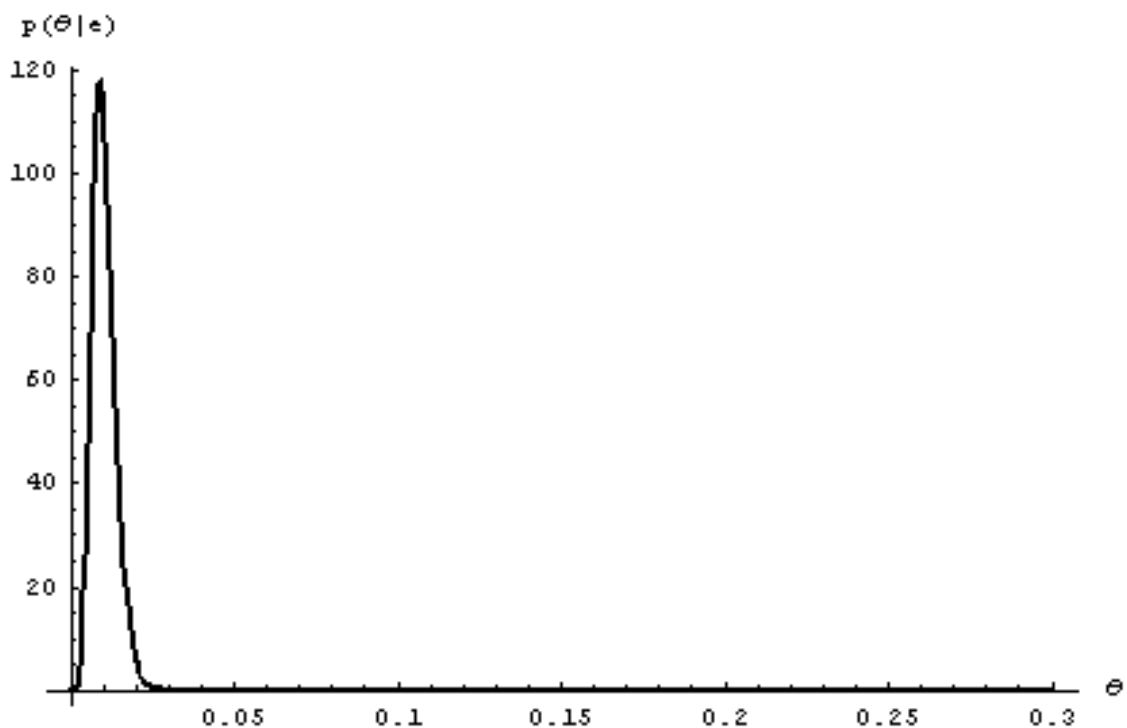


Figure 3: Expert information

It is apparent that there is virtually no probability on the long right tail. A closer view of the relevant part of the prior is graphed in Figure 4.

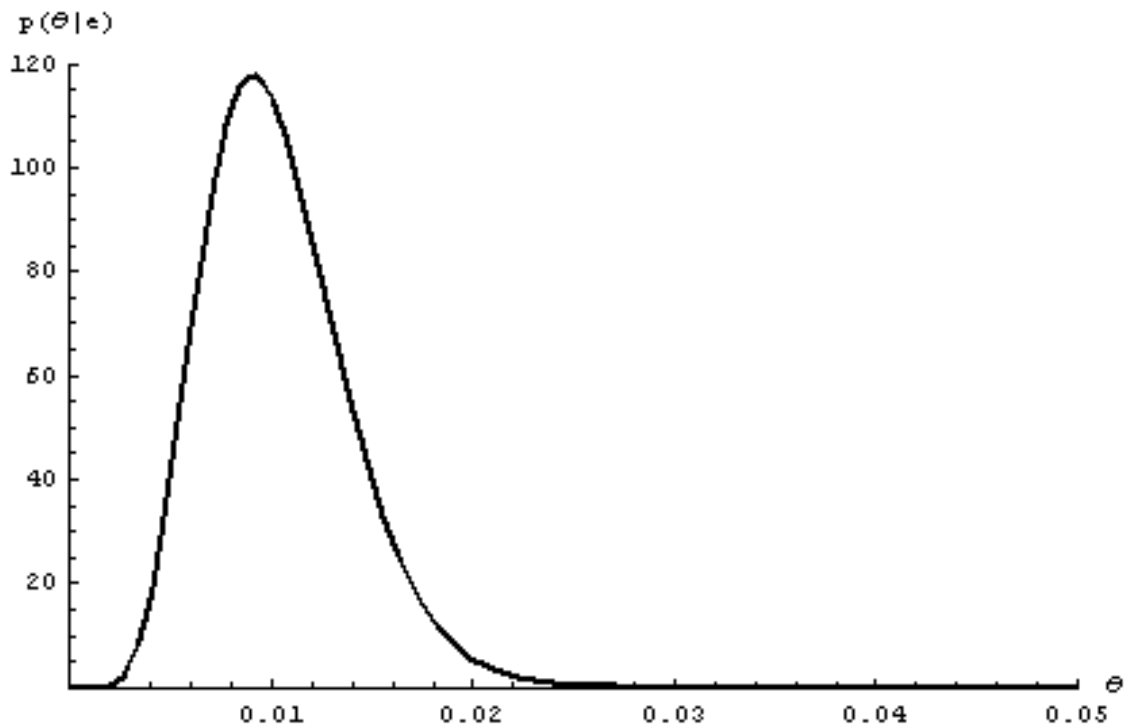


Figure 4: Expert information (closeup)

The median of this distribution is 0.00988, the mean is 0.0103 and the standard deviation is 0.00355. In practice, after the information is aggregated into an estimated probability distribution, then additional properties of the distribution would be calculated and the expert would be consulted again to see if any changes were in order before proceeding to data analysis (Lindley 1982a). This process would be repeated as necessary. In the present application there was one round of feedback, valuable since the expert had had time to consider the probabilities involved. The characteristics reported are from the second round of elicitation. An application within a bank might do additional rounds with the expert, or consider alternative experts and a combined prior.

The predictive distribution 5.1 corresponding to this prior is given in Figure 5 for $n=500$.

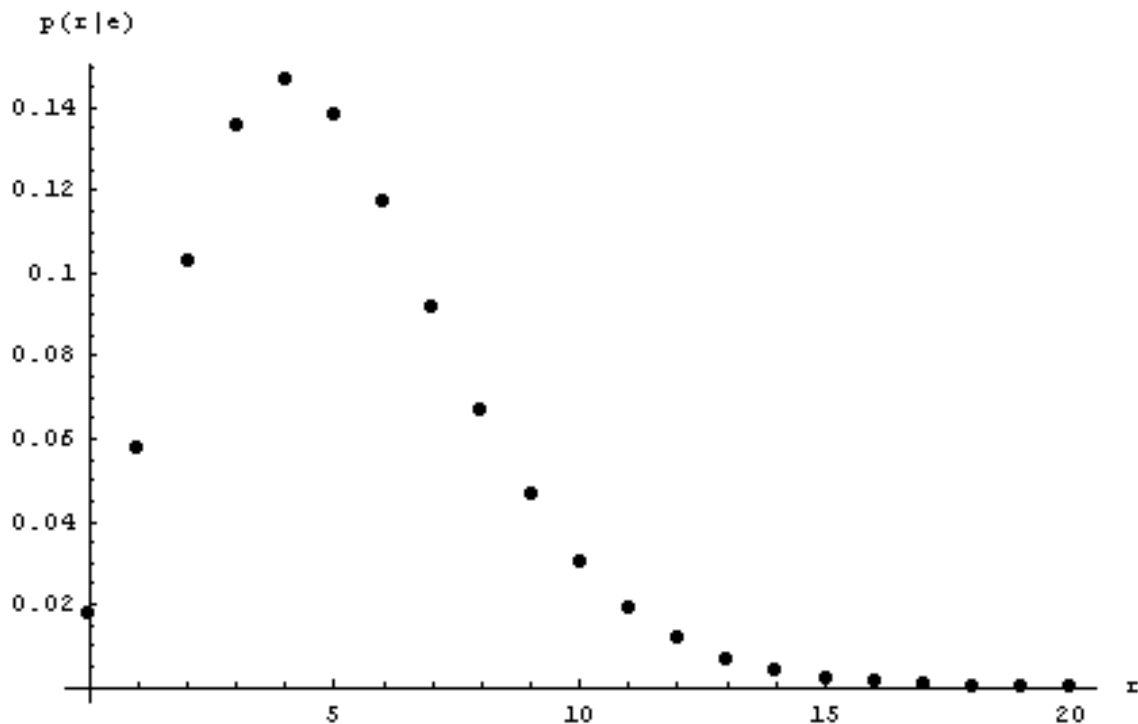


Figure 5: Predictive distribution $p(r, e)$

With our specification, the expected value of r , $E(r|e) = \sum_{k=0}^n kp(k|e)$ is 5.1 for $n=500$. Total defaults numbering 0-9 characterize 92% of expected data sets. Thus, carrying out our analysis for these 10 data types, comprising about 2^{62} distinct datasets, a trivial fraction of the 2^{500} possible datasets, actually covers 92% of the expected realizations. Defaults are expected to be rare events. This is the key to the ALD approach: we are not analyzing one particular dataset, rather we provide results applicable to 92% of the likely datasets.

7 Posterior Analysis

The posterior distribution, $p(\theta|r, e)$, is graphed in Figure 6 for $r = 0, 2, 4, 6,$ and 8 and $n=500$. The corresponding likelihood functions, for comparison, were given in figures 1 and 2. Note the substantial differences in location. Comparison with the likelihood functions graphed in Figure 1 and the prior distribution graphed in Figure 3 reveals that the expert provides much more information to the analysis than do the data.

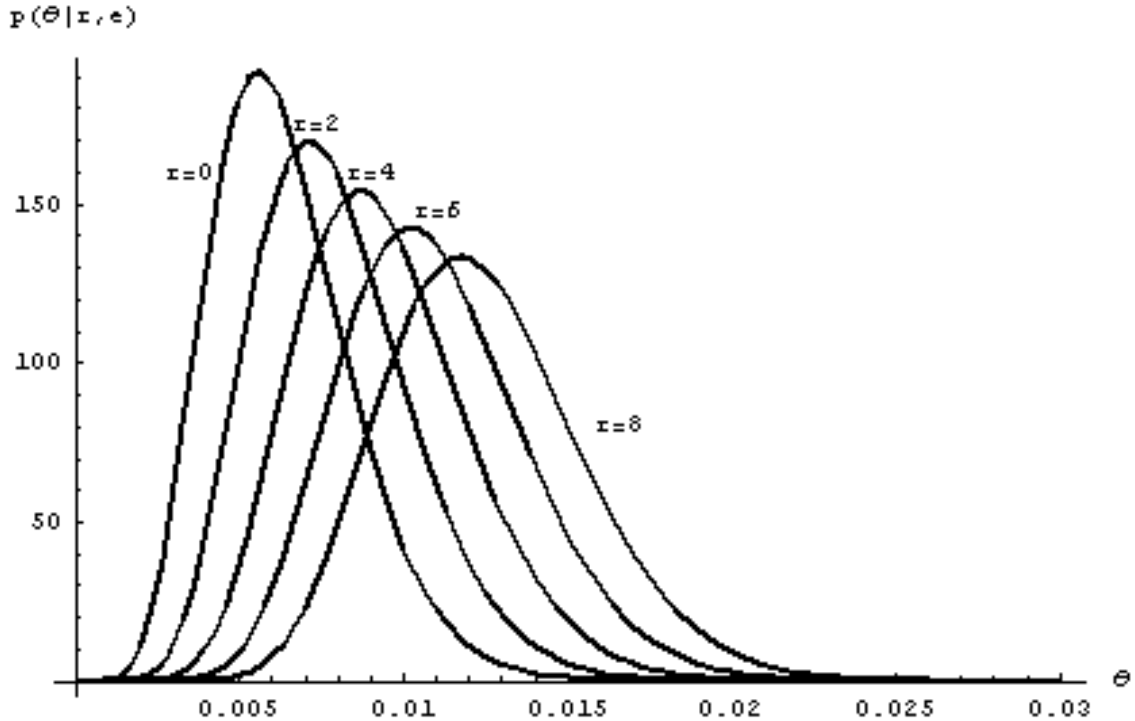


Figure 6: Posterior densities for typical samples

Given the distribution $p(\theta|r, e)$, we might ask for a summary statistic, a suitable estimator for plugging into the required capital formulas as envisioned by (Basel Committee on Banking Supervision 2004). A natural value to use is the posterior expectation, $\bar{\theta} = E(\theta|r, e)$. The expectation is an optimal estimator under quadratic loss and is asymptotically an optimal estimator under a wide variety of loss functions. An alternative, by analogy with the maximum likelihood estimator $\hat{\theta}$, is the posterior mode θ . As a summary measure of our confidence we would use the posterior standard deviation $\sigma_{\theta} = \sqrt{E(\theta - \bar{\theta})^2}$. By comparison, the usual approximation to the standard deviation of the maximum likelihood estimator is $\sigma_{\hat{\theta}} = \sqrt{\hat{\theta}(1 - \hat{\theta})/n}$. These quantities are given in Table 1 for $r=0-9$ and $r=20, 50, 100, 200$. As noted, the $r=0-9$ case covers the 2^{62} most likely datasets out of the possible 2^{500} . Together, these comprise analyses of 92% of likely datasets. The $r=20$ case is an extremely low probability outcome - less than 0.0001 - and is included to show the results in this case. There are approximately 2^{118} datasets corresponding to $r=20$. The rows for $r=50, 100,$ and 200 are included as a further "stress test" and will be discussed below. Their combined prior probability of occurrence is less than 10^{-14} .

n	r	$\bar{\theta}$	$\dot{\theta}$	$\hat{\theta}$	σ_{θ}	$\sigma_{\hat{\theta}}$
500	0	0.0063	0.0081	0.000	0.0022	0 (!).
500	1	0.0071	0.0092	0.002	0.0023	0.0020
500	2	0.0079	0.0103	0.004	0.0025	0.0028
500	3	0.0086	0.0114	0.006	0.0026	0.0035
500	4	0.0094	0.0125	0.008	0.0027	0.0040
500	5	0.0102	0.0136	0.010	0.0028	0.0044
500	6	0.0109	0.0147	0.012	0.0029	0.0049
500	7	0.0117	0.0158	0.014	0.0030	0.0053
500	8	0.0125	0.0169	0.016	0.0031	0.0056
500	9	0.0132	0.0180	0.018	0.0032	0.0060
500	20	0.0215	0.0296	0.040	0.0040	0.0088
500	50	0.0431	0.0425	0.100	0.0053	0.0134
500	100	0.0753	0.0749	0.200	0.0065	0.0179
500	200	0.1267	0.1266	0.400	0.0069	0.0219

Table 1: Default Probabilities: Location and Precision

For values of r below its expected value the posterior mean is greater than the MLE, for values above the posterior is less than the MLE, as expected. As is well-known and widely discussed, the MLE is unsatisfactory when there are no observed defaults (Basel Committee on Banking Supervision 2005), (Pluto and Tasche 2005), (BBA, LIBA, and ISDA 2005), (Kiefer 2006a). The Bayesian approach provides a coherent resolution of the inference problem without resort to desperation (sudden reclassification of defaulted assets, technical gimmicks).

Expert information will have larger weight in smaller sample sizes, and smaller relative weight for larger sample sizes. For $n=1000$, for example, $r=5-15$ reflects 76% of the most likely datasets; $r=0-20$ represents 97%. To put this in perspective, the cases $r=0-20$ correspond to approximately 2^{138} datasets out of a possible 2^{1000} . Thus, 97% of the likely observations are contained in the small fraction 2^{-862} of the possible datasets, or 0.0021 of the possible types. A substantial simplification results from concentrating on the distribution of the sufficient statistic and use of expert judgement to characterize possible samples. Naturally, this simplification depends critically on the use of expert judgement in specification of the likelihood function (our choice admits a sufficient statistic) and in specification of the prior

distribution. Rather than resorting to extensive tabulation, we report ALD results for 97% of likely samples in Figure 7. The error bands, dotted for the MLE and dashed for the prior mean, are plus/minus one standard deviation.

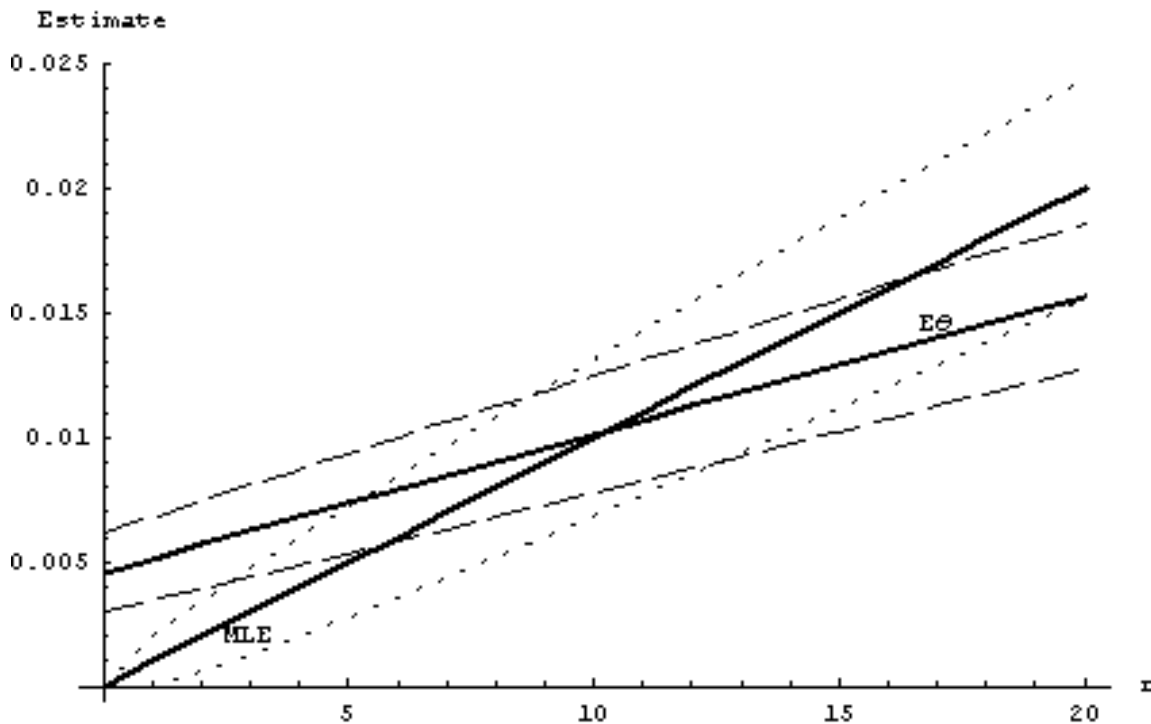


Figure 7: $E\theta$ and MLE for $n=1000$

Turning now to an extremely large sample, in which inference is not quite so problematic, as the likelihood can be expected to dominate the prior, we find a lessened role for the expert. With $n=10000$, $r=45-155$ covers 88% of all datasets. In these cases the likelihood and Bayesian analyses essentially coincide. Estimators and associated error bands for the 88% ALD analysis are shown in Table 8.

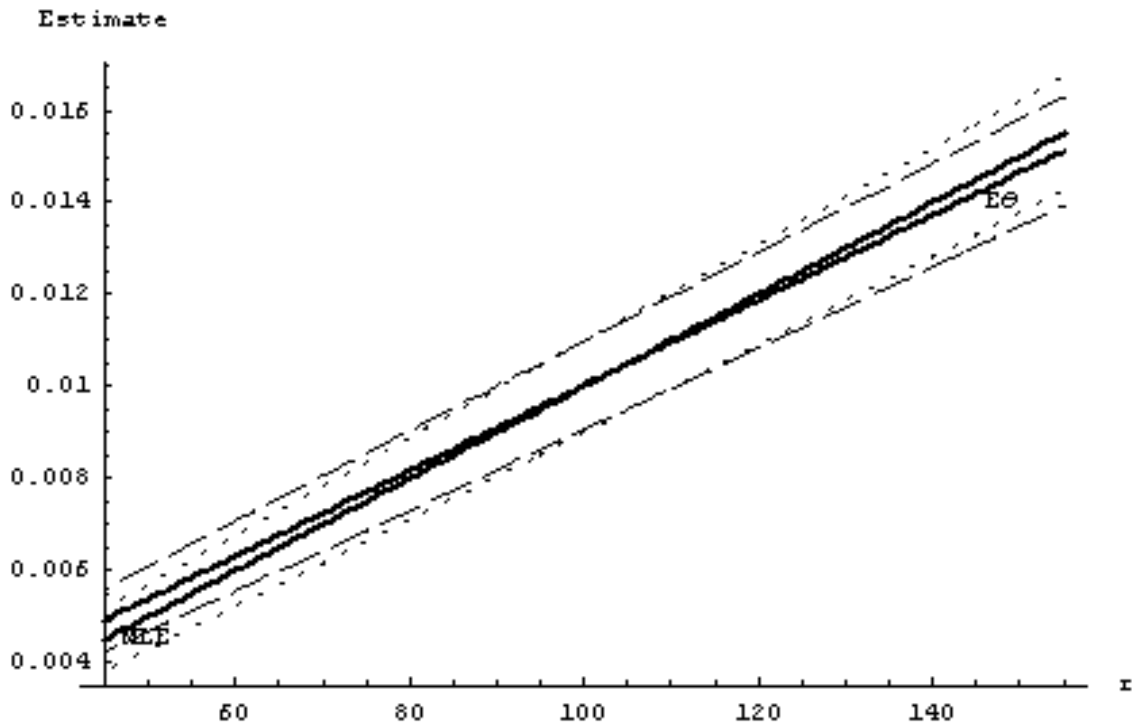


Figure 8: $E\theta$ and MLE for $n=10000$

There is still a clear difference for extremely unlikely values of r ; thus for $r=0$, $E\theta = 0.00083$, while the MLE is zero.

8 Robustness - the cautious Bayesian

Suppose we are rather less sure of our expert than he is of the default probability. Or, more politely, how can we assess just how important the tightly-held views of the expert are in determining our estimates? Of course, Table 1 gives an answer by comparing the MLE and the posterior location measures. Another answer was proposed by (Kiefer 2006b), who considered a less-certain expert with a prior with the same location but substantially higher variance than the actual expert. An alternative approach, rather more formal and based on the literature on Bayesian robustness (Berger and Berliner 1986) is to mix the actual expert's prior with an alternative prior, and see exactly how seriously the inferences are affected by changes in the mixing parameter. (Berger and Berliner 1986) in fact suggesting mixing in a class of distributions, corresponding to different amounts or directions of uncertainty in the prior elicitation. In this spirit, we will mix the expert's 4-parameter

beta distribution with a uniform distribution. Here, there are two clear possibilities. One is to mix with the uniform on $[a,b]$, accepting the expert's bounds but examining robustness to alpha and beta. The second is to mix with the uniform on $[0,1]$, allowing all theoretically feasible values of θ . We choose the latter approach. This is not a completely comfortable approach. Although the uniform is commonly interpreted as an uninformative prior, it in fact has a mean of $1/2$, not a likely value for our default probability by any reasonable prior. An alternative might be to mix with a prior with the same mean as our expert's distribution, but maximum variance. We do not pursue this here. Our results suggest that it would not make much difference; the key is to mix in a distribution with full support, so that likelihood surprises can appear. We choose to mix the expert's prior with a uniform on all of $[0,1]$. This allows input from the likelihood if the likelihood happens to be concentrated above b (or below a). The mixture distribution is

$$p(\theta|e, \epsilon) = (1 - \epsilon)p(\theta|\alpha, \beta, a, b)I(\theta \in [a, b]) + \epsilon \quad (8.1)$$

for $\theta \in [0, 1]$. Here $I(\theta \in [a, b])$ is the indicator function equal to one if the argument is true, zero otherwise. The approach can be used whatever prior is specified, not just the 4-parameter beta. Our robust prior is in the 9-parameter mixture family 4.3, consisting of our expert's 4-parameter beta mixed with the 4-parameter beta with parameters $\{\alpha, \beta, a, b\} = \{1, 1, 0, 1\}$ and mixing parameter ϵ . Table 2 shows the posterior means for the mixture priors for $\epsilon = \{0.01, 0.1, 0.2, 0.3, 0.4\}$.

n	r	$\bar{\theta}; \epsilon = .01$	$\bar{\theta}; \epsilon = .1$	$\bar{\theta}; \epsilon = .2$	$\bar{\theta}; \epsilon = .3$	$\bar{\theta}; \epsilon = .4$
500	0	0.0063	0.0063	0.0062	0.0061	0.0061
500	1	0.0071	0.0071	0.0071	0.0071	0.0070
500	2	0.0079	0.0079	0.0079	0.0079	0.0078
500	3	0.0086	0.0086	0.0086	0.0086	0.0086
500	4	0.0094	0.0094	0.0094	0.0094	0.0094
500	5	0.0102	0.0102	0.0102	0.0102	0.0102
500	6	0.0109	0.0109	0.0110	0.0110	0.0110
500	7	0.0117	0.0117	0.0117	0.0118	0.0118
500	8	0.0125	0.0125	0.0125	0.0125	0.0126
500	9	0.0132	0.0133	0.0133	0.0134	0.0134
500	20	0.0358	0.0358	0.0386	0.0398	0.0405
500	50	0.1016	0.1016	0.1016	0.1016	0.1016
500	100	0.2012	0.2012	0.2012	0.2012	0.2012
500	200	0.4004	0.4004	0.4004	0.4004	0.4004

Mixing in the uniform prior makes essentially no difference to the posterior mean for data in the likely part of the set of potential samples. For $r=20$, unlikely but not outrageous, mixing in the prior makes a substantial difference. For the extremely unlikely values, 50, 100, 200, the differences are dramatic. The actual value of ϵ makes almost no difference. The numbers for $\epsilon = 0.001$, not shown in the table, give virtually the same mean for all r . For comparison, we recall the values of $\bar{\theta}$ for $r=\{20,50,100,200\}$ from Table 1. These are $\{0.0215,0.0431,0.0753,0.1267\}$. Figure 9 shows the posterior distributions for our expert's prior, $p(\theta|r, e)$ for $r=50, 100$, and 200. It is clear that the prior plays a huge role here, as the likelihood mass is concentrated near .1, .2 and .4, while the prior gives only trivial weight to values greater than about .03, see Figures 1 and 3. On the other hand, Figure 10 shows the posterior corresponding to 8.1 with 1% mixing ($\epsilon = 0.01$). Here, the likelihood dominates, as the likelihood value near the expert's prior is vanishingly small relative to the likelihood in the tail area of the mixing prior.

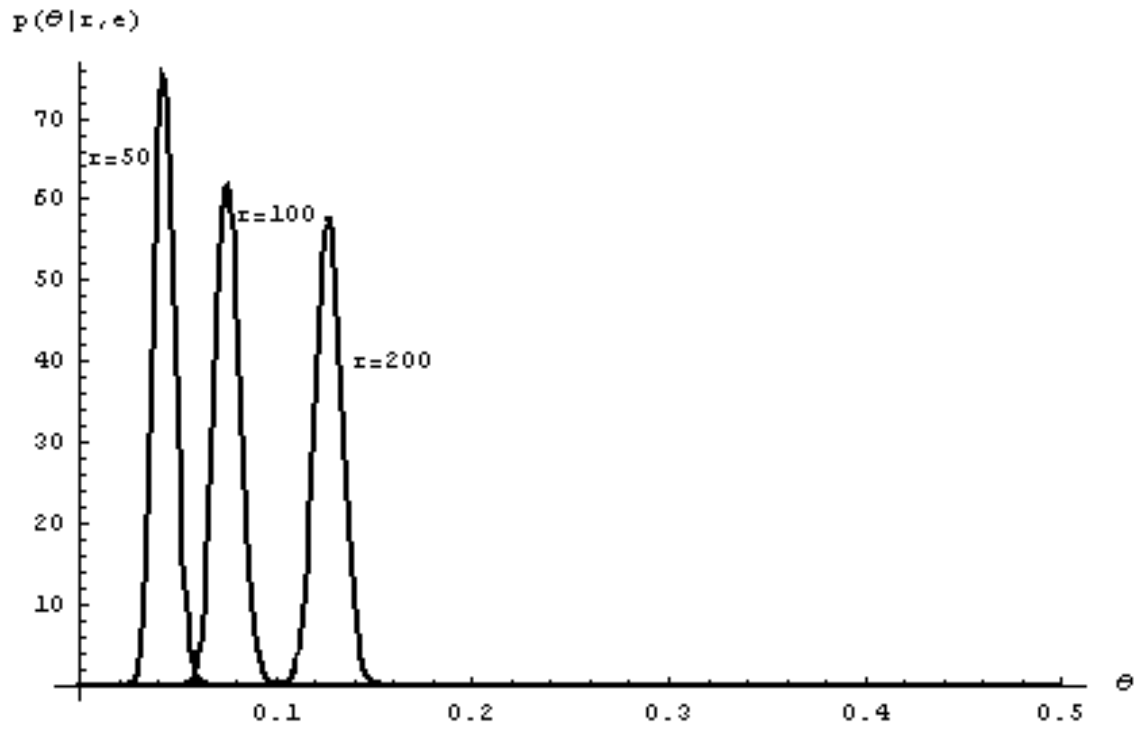


Figure 9: Posteriors using 8.1 with $\epsilon = 0$ and $r=50,100$, and 200

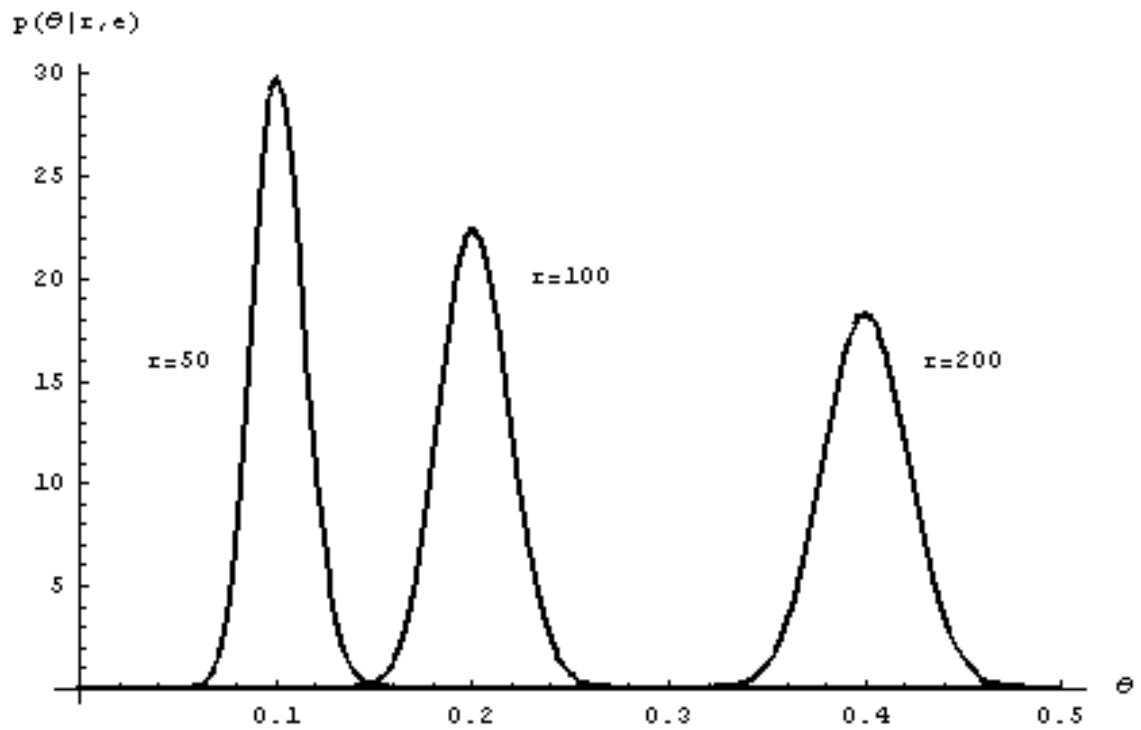


Figure 10: Posteriors using 8.1 with $\epsilon = 0.01$ and $r = 50, 100,$ and 200 .

Thus, the robust analysis with even a very small nonzero mixing fraction can reveal disagreements between the data and the expert opinion which are perhaps masked by the formal analysis. This robust analysis may have a role to play in the validation phase.

I what sense is the robust analysis useful? We are really bringing something outside the model, namely the uniform distribution representing no one's beliefs, into the analysis as a formal tool for diagnostic analysis. The spirit is the same as usual procedures associated with good statistical practice - residual analysis, out of sample fits, forecast monitoring, or comparison with alternative models. All of these procedures involve stepping away from the specified model and its analysis, and asking, post estimation, does the specification make sense? Can we tell where it breaks down? Can it be improved? These post-estimation model evaluation techniques are often informal, sometimes problem specific, and require sound statistical judgement (OCC 2006). The analysis of robustness via an artificial prior is an attempt to merge the formal analysis with the informal post-estimation model checking. A related method, checking for irrelevant data using a mixture distribution, is proposed by (Ritov 1985) and this might have a role as well.

9 Conclusion

I have considered inference about the default probability for a midrange portfolio segment on the basis of data information and expert judgement. Examples focus on the sample size of 500; results are also presented for the large sample sizes of 1000 and 10000 observations, not unreasonable for large banks in this risk range. These analyses are relevant to hypothetical portfolios of middle-risk commercial loans. These are predominantly to unrated companies; if rated these would be approximately S&P Baa or Moody's BBB. I have also represented the judgement of an expert in the form of a probability distribution, for combination with the likelihood function. The expert is a practitioner experienced in risk management in well-run banks. The 4-parameter Beta distribution seems to reflect expert opinion fairly well. Errors, which would be corrected through additional feedback and re-specification in practice, are likely to introduce more certainty into the distribution

rather than less. There are no real data here; the portfolios are hypothetical. However, using the ALD approach, it is possible to study the posterior distributions for all of the most likely configurations of defaults in the samples. Using ALD, we consider the possible realizations of the sufficient statistic for the specified statistical model. In the default case, the number of realizations is linear in the sample size (while of course the number of potential distinct samples is exponential). Using the expert information, it is possible to isolate the most likely realizations. In the sample of 500, five defaults are expected. In this case, our analysis of 0 through 9 defaults covers 92% of expected datasets. Our analyses of samples of 1000 and 10000 covered 97% and 88% of the likely realizations respectively.

At the validation stage, modelers can be expected to have to justify the likelihood specification and the representation of expert information. Analysis of the sensitivity of the results to the prior should be a part of this validation procedure. We propose using a mixture of the expert's prior and an alternative, less informative prior. In our case, we mix the prior with a uniform distribution on the unit interval. While it is not likely that the uniform describes any expert's opinion on the default probability, mixing in the uniform allows unexpected disagreement between the prior and the data to appear vividly. An example shows that even a trivially small weight on the alternative will do. Of course, within the context of the model, the decision based on the expert's posterior is correct. A broader view might suggest something wrong with the specification - of either the likelihood or the prior. Perhaps these do not refer to the same risk class, or perhaps the default definitions are inconsistent. The situation is not unlike that arising in ordinary validation exercises in which the model is evaluated in terms of residual analysis or out-of-sample fits. These involve considerations which are relevant but which are outside the formal model. As a result there are a number of different methods in use, corresponding to different ways in which models can fail, and expert judgement remains crucial in this less formal context as well as in the formal specification of the likelihood and the prior (OCC 2006).

References

BASEL COMMITTEE ON BANKING SUPERVISION (2004): "International Convergence of Capital Measurement and Capital Standards: A Revised Framework,"

- Bank for International Settlements.
- (2005): “Basel Committee Newsletter No. 6: Validation of Low-Default Portfolios in the Basel II Framework,” Discussion paper, Bank for International Settlements.
- BBA, LIBA, AND ISDA (2005): “Low Default Portfolios,” Discussion paper, British Banking Association, London Investment Banking Association and International Swaps and Derivatives Association, Joint Industry Working Group.
- BERGER, J., AND L. M. BERLINER (1986): “Robust Bayes and Empirical Bayes Analysis with Contaminated Priors,” *The Annals of Statistics*, 14(2), 461–486.
- CHALONER, K. M., AND G. T. DUNCAN (1983): “Assessment of a Beta Prior Distribution: PM Elicitation,” *The Statistician*, 32(1/2, Proceedings of the 1982 I.O.S. Annual Conference on Practical Bayesian Statistics), 174–180.
- COVER, T. M., AND J. A. THOMAS (1991): *Elements of Information Theory*. John Wiley & Sons.
- DE FINETTI, B. (1974): *Theory of Probability*, vol. 1. New York: Wiley.
- DIACONIS, P., AND D. YLVISAKER (1985): “Quantifying Prior Opinion,” in *Bayesian Statistics 2*, ed. by J. M. Bernardo, M. H. DeGroot, D. Lindley, and A. Smith, pp. 133–156. Elsevier Science Publishers BV (North-Holland).
- GARTHWAITE, P. H., J. B. KADANE, AND A. O’HAGAN (2005): “Statistical Methods for Eliciting Probability Distributions,” *Journal of the American Statistical Association*, 100, 780–700.
- KADANE, J. B., J. M. DICKEY, R. L. WINKLER, W. S. SMITH, AND S. C. PETERS (1980): “Interactive Elicitation of Opinion for a Normal Linear Model,” *Journal of the American Statistical Association*, 75(372), 845–854.
- KADANE, J. B., AND L. J. WOLFSON (1998): “Experiences in Elicitation,” *The Statistician*, 47(1), 3–19.
- KIEFER, N. M. (2006a): “Default Estimation for Low Default Portfolios,” OCC Working Paper.

- KIEFER, N. M. (2006b): “The Probability Approach to Default Estimation,” Discussion paper, Cornell University.
- LINDLEY, D. V. (1982a): “The Improvement of Probability Judgements,” *Journal of the Royal Statistical Society. Series A (General)*, 145(1), 117–126.
- LINDLEY, D. V. (1982b): “Scoring Rules and the Inevitability of Probability,” *International Statistical Review / Revue Internationale de Statistique*, 50(1), 1–11.
- OCC (2006): “Validation of Credit Rating and Scoring Models: A workshop for Managers and Practitioners,” in *Validation of Credit Rating and Scoring Models*.
- PLUTO, K., AND D. TASCHE (2005): “Thinking Positively,” *Risk*, August, 72–78.
- RITOV, Y. (1985): “Robust Bayes Decision Procedures: Gross Error in the Data Distribution,” *The Annals of Statistics*, 13(2), 626–637.