

# Risk Based Scores and the Gini Index

Edward W. (Jed) Frees <sup>\*</sup>      Glenn Meyers<sup>†</sup>      A. David Cummings<sup>‡§</sup>

April 27, 2011

*Abstract.* In 1905, Max Otto Lorenz displayed skewed income distributions using a graph now known as the Lorenz curve. In 1912, Corrado Gini summarized this curve with a statistic now known as the Gini index. Both devices are widely used in welfare economics, among other fields. In this paper, we extend these concepts to a financial context by ordering risks using relativities which are risk based scores relative to prices.

Using the relativity ordering, we develop a Lorenz curve and the corresponding Gini index that can cope with adverse selection and measure potential profit. We provide a detailed example using personal lines homeowners insurance. Further, we show that the Gini index can be written in terms of covariance operators, thus expanding the scope of interpretations. We implement theory developed in a companion paper to calibrate sample sizes, establishing that the number of observations typically encountered in insurance practice are sufficient for reliable estimation of our new Gini index.

---

<sup>\*</sup>University of Wisconsin and ISO Innovative Analytics

<sup>†</sup>ISO Innovative Analytics

<sup>‡</sup>ISO Innovative Analytics

<sup>§</sup>Keywords: Insurance pricing, association measures, adverse selection, Lorenz curve

# 1 Introduction

We wish to compare the distribution of a financial risk  $y$  to a price  $P$ . We assume that the analyst has available a set of known exogenous risk characteristics  $\mathbf{x}$  upon which both the risk and price distributions depend. Our goal is to develop a measure that quantifies the extent to which  $P$  can be used to assess the distribution of  $y$ . With such a measure, we could compare alternative pricing structures. Although this modeling framework is applicable broadly, for concreteness we focus on an insurance loss as the financial risk and the price is a premium that a policyholder would pay for an insurer to cover the risk. In this special case, characteristics  $\mathbf{x}$  typically influence both distributions. Insurance scores, used for setting premiums, that are based on knowledge of  $\mathbf{x}$  are known as “risk based scores.”

**Measures of Association.** Classic statistics offers analysts many measures that quantify the association between  $y$  and  $P$ . For example, it is common to use a Pearson (the usual, or product-moment) correlation or a Spearman correlation (that quantifies correlations between ranks of variables) to assess association. Some alternative measures of association are based on quantifying the distribution of the difference between losses and premiums. For example, if one were interested in the bias of a premium, one could look to the mean absolute error statistic given as  $MAE = n^{-1} \sum_{i=1}^n |y_i - P_i|$ . An alternative is to look through mean square errors, such as through the root mean square error statistic  $RMSE = \sqrt{n^{-1} \sum_{i=1}^n (y_i - P_i)^2}$ . Of course, many variations of these basic statistics are available; for example, it is common to examine ratios in lieu of differences.

Although useful, these classic statistics are limited in our financial context for at least three reasons. First, their optimality properties are motivated by theory based on symmetric distributions (such as the normal and LaPlace distributions) and do not account for the complexity of risk distributions. To illustrate this complexity, in typical homeowners insurance data described in Section 3, 94% of the losses are zeros (corresponding to no claims) and when losses are positive, the distribution tends to be right-skewed and thick-tailed. Second, they do not explicitly allow for the asymmetric nature of decisions faced by analysts. For example, one would only introduce a new premium methodology into an existing rating structure if the new methodology were dramatically better than an existing alternative. Third, these classic statistics lack economic interpretation.

In this paper, we develop summary measures that can address these limitations. These measures are natural extensions of the classic Lorenz curve and associated Gini index, given as follows.

## 1.1 The Lorenz Curve and Gini Index

In welfare economics, it is common to compare distributions via the *Lorenz curve*, developed by Max Otto Lorenz (1905). A Lorenz curve is a graph of the proportion of a population on the horizontal axis and a distribution function of interest on the vertical axis. It is typically used to represent income distributions. When the income distribution is perfectly aligned with the population distribution, the Lorenz curve results in a 45 degree line that is known as the “line of equality.” The area between the Lorenz curve and the line of equality is a measure of the discrepancy between the income and population distributions. Two times this area is known as the *Gini index*, introduced by Corrado Gini in 1912. See, for example, Sen and Foster (1998), for additional background on income equality. For readers interested in examining current international inequality measures, see the online resource UNU-Wider World Income Inequality Database (2008).

The contributions of Joseph Gastwirth in the 1970’s (e.g., Gastwirth, 1971, 1972) helped to emphasize the importance of the Lorenz curve and the Gini index as tools for comparing distributions, particularly in economic applications. The subsequent literature is extensive. In one strand of the literature, researchers have sought to understand differences in economic equality among population subgroups (e.g., Lambert and Decoster, 2005, Gastwirth, 1975). In another strand, analysts have introduced weight functions into the Lorenz curve (e.g., to account for the number of publications when studying impact factors, Egghe, 2005). Yitzhaki (1996) describes how weighted regression sampling estimators can be of interest in welfare economic applications. Here, the idea is to adjust regression weights for social attitudes toward inequality. In another stream of research, analysts have used the Gini index for model selection in genomics (Nicodemus and Malley, 2009) and in classification trees (Sandri and Zuccolotto, 2008).

**Example: Distribution of Homeowners Premiums.** For an insurance example, Figure 1 shows a distribution of premiums. This figure is based on a sample of 359,454 policyholders with premiums that will be described in Section 3. The left-hand panel shows a right-skewed histogram of premiums. When plotting this figure, premiums that exceeded 1,200 were ignored. The right-hand panel provides the corresponding Lorenz curve, showing again a skewed distribution. For example, the arrow marks the point where 60% of the policyholders pay 40% of premiums. The 45 degree line is the “line of equality;” if each policyholder paid the same premium, then the premium distribution would be at this line. The Gini index, twice the area between the Lorenz curve and the 45 degree line, is 29.5% for this data set.

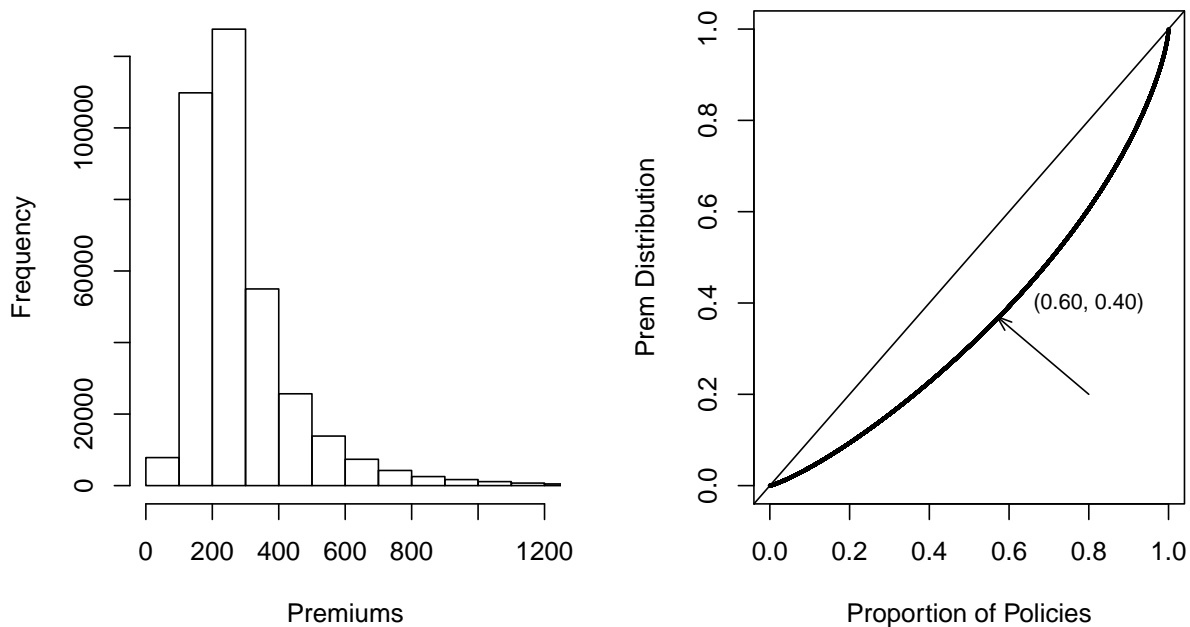


Figure 1: Distribution of Premiums. The left-hand panel is a histogram of premiums from a group of 359,454 policyholders, showing a distribution that is right-skewed. The right-hand panel provides the corresponding Lorenz curve. The arrow marks the point where 60% of the policyholders pay 40% of premiums.

## 1.2 Relating Premium to Loss Distributions

From Section 1.1, the Lorenz curve is a device that is well-known in welfare economics for displaying distributions. It is particularly useful for interpreting skewed distributions, a shape that insurance analysts are well acquainted with.

One could use classic Lorenz curves to compare a premium to a loss distribution. For example, it would be straightforward to compute Lorenz curves for premiums and for losses, and then superimpose the two curves on the same figure. However, the population distribution for each curve would be based on different sort orders (by premiums and losses, respectively), so that it would not be meaningful to compare premiums to losses for any policyholder group.

**The Role of Relativities.** As an alternative, in the following section we extend the Lorenz curve through the introduction of a third variable called a *relativity*. The relativity connects the losses

to the premiums and is the variable that we will sort on, thus maintaining consistency between policyholder groups. In this way, we can track differences between losses and premiums and, through different sort orders, emphasize the differences between these two distributions. Because premiums (but not losses) can be influenced by insurance analysts, we will argue that this comparison provides a way to judge whether a given premium  $P$  is somehow “better” than an alternative.

The plan for the paper follows. In Section 2 we develop this extension of the Lorenz curve and Gini index, providing definitions, giving an example and focusing on a special case of interest. We developed theoretical properties of this Gini index elsewhere (Frees, Meyers and Cummings, 2011a); to keep this paper self-contained, these properties are summarized in Appendix Section 8. Section 3 provides a detailed example using a sample of policies from homeowners insurance that shows how one can use the new Gini index. Section 4 provides additional interpretations of the Gini index, showing how it can be expressed in terms of covariance functions. Details of the covariance calculations are in Appendix Section 9. Section 5 summarizes a small simulation study that shows how to use the Gini statistic as a tool for model selection. Section 6 then describes how one can use the Gini index to suggest an appropriate sample size for pilot testing, with supporting calculations in Appendix Section 10. Section 7 closes with a summary and some additional remarks.

## 2 Ordered Lorenz Curve and the Gini Index

We now introduce an *ordered* Lorenz curve which is a graph of the distribution of losses versus premiums, where both losses and premiums are ordered by relativities. Intuitively, the relativities point towards aspects of the comparison where there is a mismatch between losses and premiums. To make the ideas concrete, we first provide some notation. We will consider  $i = 1, \dots, n$  policies. For the  $i$ th policy, let

- $y_i$  denote the insurance loss,
- $\mathbf{x}_i$  be the set of policyholder characteristics known to the analyst,
- $P_i = P(\mathbf{x}_i)$  be the associated premium that is a function of  $\mathbf{x}_i$ ,
- $S_i = S(\mathbf{x}_i)$  be an insurance score under consideration for rate changes, and
- $R_i = R(\mathbf{x}_i) = S(\mathbf{x}_i)/P(\mathbf{x}_i)$  is the “relativity,” or relative premium.

Thus, the set of information used to calculate the ordered Lorenz curve for the  $i$ th policy is  $(y_i, P_i, S_i, R_i)$ .

**Ordered Lorenz Curve.** We now sort the set of policies based on relativities (from smallest to largest) and compute the premium and loss distributions. Using notation, the premium distribution is

$$\hat{F}_P(s) = \frac{\sum_{i=1}^n P(\mathbf{x}_i) I(R_i \leq s)}{\sum_{i=1}^n P(\mathbf{x}_i)}, \quad (1)$$

and the loss distribution is

$$\hat{F}_L(s) = \frac{\sum_{i=1}^n y_i I(R_i \leq s)}{\sum_{i=1}^n y_i}, \quad (2)$$

where  $I(\cdot)$  is the indicator function, returning a 1 if the event is true and zero otherwise. The graph  $(\hat{F}_P(s), \hat{F}_L(s))$  is an *ordered Lorenz curve*.

**Example: Homeowners Loss and Premium Distributions.** As an example of an ordered Lorenz curve, Figure 2 shows a curve using our homeowners data. For this curve, the score “SP\_FreqSev\_Basic” was used as the base premium and the score “IND\_FreqSev” was used to compute the relativities. To help interpret the curve, an arrow marks a typical point, corresponding to 60% of premium and 53.8% of losses. That is, with knowledge of relativities, an insurer could identify a portfolio that enjoys 60% of premiums with only 53.8% of losses. This is a profitable portfolio, one well worth retaining.  $\square$

**The Role of Adverse Selection.** If an insurer does not adopt a refined rating plan but its competitors do, then the insurer could become victim of adverse selection. The insurer’s good risks will switch to the competitor, and the insurer’s remaining risks will have higher losses. The ordered Lorenz curve quantifies the extent of this potential adverse selection. Through ordering of the relativity, it summarizes the performance of portfolios that are profitable and hence subject to potential raiding by competitors. Conversely, it also summarizes the performance of poorly performing portfolios that an insurer may wish to examine for potential loss control activities, e.g., tighter underwriting restrictions.

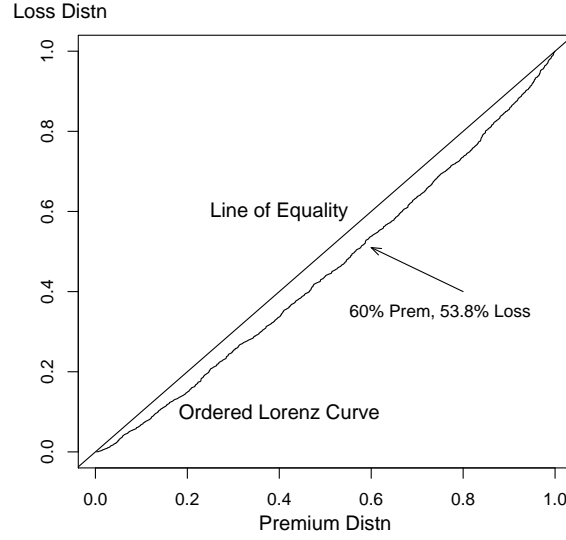


Figure 2: An Ordered Lorenz Curve. For this curve, the corresponding Gini index is 10.03% with a standard error of 1.45% .

**The Gini Index.** Of course, the selection of the 60th premium percentile is arbitrary. Insurers will wish to consider the profitability of different size portfolios. Thus, we summarize the curve using the Gini index which is (twice) the area between the curve and the 45 degree line, known as “the line of equality.” The line of equality can be interpreted as a “break-even” case for the insurer, where the percentage of losses equals the percentage of premiums. Curves below the line of equality represent a profitable situation for the insurer.

Specifically, the Gini index can be calculated as follows. Suppose that the empirical ordered Lorenz curve is given by  $\{(a_0 = 0, b_0 = 0), (a_1, b_1), \dots, (a_n = 1, b_n = 1)\}$  for a sample of  $n$  observations. Here, we use  $a_j = \hat{F}_P(R_j)$  and  $b_j = \hat{F}_L(R_j)$ . Then, the empirical Gini index is

$$\begin{aligned} \widehat{Gini} &= 2 \sum_{j=0}^{n-1} (a_{j+1} - a_j) \left\{ \frac{a_{j+1} + a_j}{2} - \frac{b_{j+1} + b_j}{2} \right\} \\ &= 1 - \sum_{j=0}^{n-1} (a_{j+1} - a_j)(b_{j+1} + b_j). \end{aligned} \quad (3)$$

As described in Section 1.1, the classic Lorenz curve shows the proportion of policyholders on the horizontal axis and the loss distribution function on the vertical axis. The “ordered” Lorenz curve extends the classical Lorenz curve in two ways, (1) through the ordering of risks and prices by relativities and (2) by allowing prices to vary by observation. We summarize the ordered Lorenz

curve in the same way as the classic Lorenz curve using a Gini index, defined as twice the area between the curve and a 45 degree line. The analyst seeks ordered Lorenz curves that approach passing through the southeast corner (1,0); these have greater separation between the loss and premium distributions and therefore larger Gini indices.

**Example.** Suppose we have only  $n = 5$  policyholders with experience as:

Variable	$i$	1	2	3	4	5	Sum
Loss	$y_i$	5	5	5	4	6	25
Premium	$P(\mathbf{x}_i)$	4	2	6	5	8	25
Relativity	$R(\mathbf{x}_i)$	5	4	3	2	1	

Figure 3 compares the Lorenz curve to the ordered version based on this data. The left-hand panel shows the Lorenz curve. The horizontal axis is the cumulative proportion of policyholders (0, 0.2, 0.4, and so forth) and the vertical axis is the cumulative proportion of losses (0, 4/25, 9/25, and so forth). This figure shows little separation between the distributions of losses and policyholders.

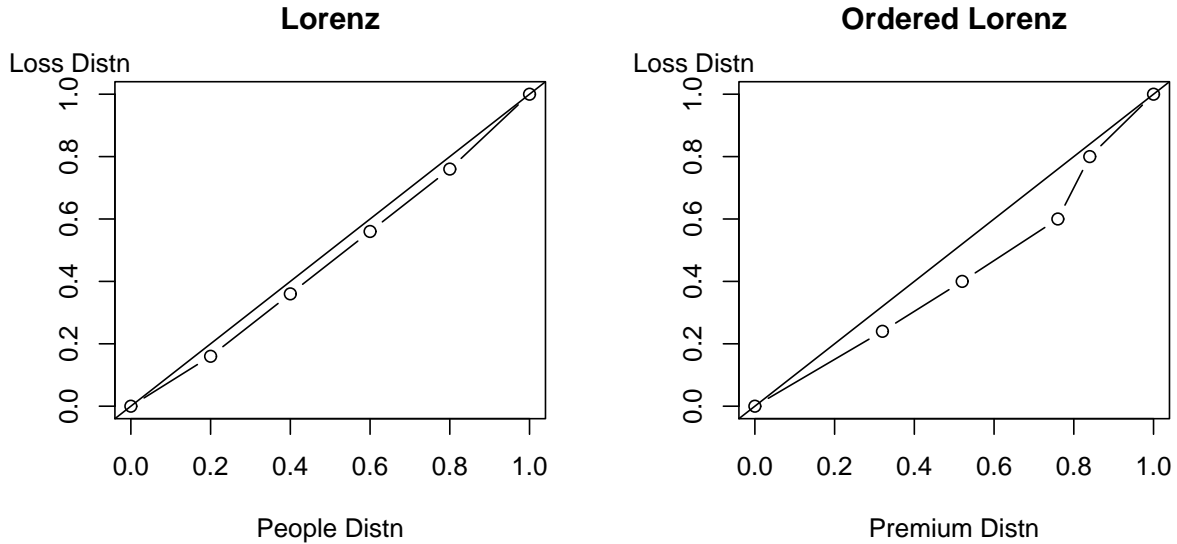


Figure 3: Lorenz versus Ordered Lorenz Curve. The Gini index for the left-hand panel is 5.6%. The Gini index for the right-hand panel is 14.9%.



The right-hand panel shows the ordered Lorenz curve. Because observations are sorted by relativities, the first point after the origin (reading from left to right) is  $(8/25, 6/25)$ . The second point is  $(13/25, 10/25)$ , with the pattern continuing. For the ordered Lorenz curve, the horizontal axis uses premium weights, the vertical axis uses loss weights, and both axes are ordered by relativities. From the figure, we see that there is greater separation between losses and premiums when viewed through this relativity.  $\square$

**Rescaling of Premiums and Losses.** From equations (1) and (2), we see that we can arbitrarily rescale premiums and losses by any positive constant and the distribution functions remain unchanged. Thus, without loss of generality, we assume henceforth that the average loss  $\bar{y}$  and average premium  $\bar{P}$  are both equal to 1.

**Properties of the Gini Index.** Appendix Section 8 describes properties the Gini index, including its consistency and asymptotic normality, that were proved in a companion paper, Frees, Meyers and Cummings (2011a). We use this asymptotic normality extensively in this paper, it is the basis for assessing the statistical significance of the Gini index.

Moreover, Frees, Meyers and Cummings (2011a) derived a result that shows that the Gini index becomes larger as one uses a “more refined” insurance score. Specifically, consider a rating plan with premiums  $P$  and an insurance score  $S$  that is determined by a regression function using a set of insured characteristics that produces a Gini index. Consider an alternative insurance score  $S_A$  that is determined by a regression function using the base set of insured characteristics plus additional information such as more precise geographic information or credit scores. Then, the Gini index produced using this refined set of information is at least as large as the Gini index using the base information.

In this sense, the Gini index provides a summary statistic that indicates whether one insurance score is better than a competitor. The paper of Frees, Meyers and Cummings (2011a) also derived estimators to judge whether Gini indices produced by alternative scoring methods are statistically different from one another.

**The Gini Index as a Measure of Profit.** One reason that the Gini index is important is because it has a direct economic interpretation. Specifically, consider an average profit,

$$\frac{1}{n} \sum_{i=1}^n \left( \hat{F}_P(R_i) - \hat{F}_L(R_i) \right) \approx \frac{\widehat{Gini}}{2}, \quad (4)$$

that can be shown to be approximately equal to the Gini index divided by two. It is an “average” in the sense that we are taking a mean over all decision-making strategies, that is, each strategy retaining the policies with relativities less than or equal to  $R_i$ . In this sense, insurers that adopt a rating structure with a large Gini index are more likely to enjoy a profitable portfolio.

Thus, we can think about the Gini index as the average expected profit to be gained by using relativities (to form portfolios). That is, the Gini index calibrates potential mismatches between losses and premiums. We change the Gini index by using different relativities; the relativities gives an idea as to where potential mismatches occur. In particular, a low relativity means that a policy is highly profitable and a good candidate to retain.

### 3 Homeowners Example

For the “gold standard” of model validation in predictive modeling, one examines the performance of a model on an independent held-out sample (e.g., Hastie, Tibshirani and Friedman, 2001). For this example, we used an in-sample dataset of 404,664 records to compute parameter estimates. We then use the estimated parameters from the in-sample model fit as well as predictor variables from a held-out, or validation, subsample of 359,454 records, whose losses we wish to predict.

#### 3.1 Comparison of Scores

More details on this database and scoring methods are available in two companion papers, Frees, Meyers and Cummings (2010, 2011b). Based on the theory developed in these two papers, we have under consideration fourteen scores that are listed in the legend of Table 2. This table summarizes the distribution of each score on the held-out data. Not surprisingly, each distribution is right-skewed.

For example, Table 2 also shows that the single-peril frequency severity model using the extended set of variables (SP\_FreqSev) provides the lowest score, both for the mean and at each percentile (below the 75th percentile). Except for this, no model seems to give a score that is consistently high or low for all percentiles. All scores have a lower average than the average held-out actual losses (TotClaims).

### 3.2 Comparing Scoring Methods to a Selected Base Premium

To compare these scoring methods, we first assume that the insurer has adopted a base premium for rating purposes; to illustrate, we use the “SP\_FreqSev\_Basic” for this premium. This method uses only a basic set of rating variables to determine insurance scores from a single-peril, frequency and severity model. Assume that the insurer wishes to investigate alternative scoring methods to understand the potential vulnerabilities of this premium base; Table 3 summarizes several comparisons using the Gini index. This table includes the comparison with the alternative score IND\_FreqSev, shown in Figure 2, as well as twelve other scores.

The standard errors are from Appendix Section 8 (derived in Frees et al., 2011a). Thus, to interpret Table 3, one may use the usual rules of thumb and reference to the standard normal distribution to assess statistical significance. For the three scores that use the basic set of variables, SP\_PurePrem\_Basic, IND\_PurePrem\_Basic, and IV\_PurePrem\_Basic, the Gini indices are between 2.5 and 4 standard errors above zero, indicating statistical significance. In contrast, the other Gini indices all are more than 7 standard errors above zero, indicating that the ordering used by each score helps detect important differences between losses and premiums.

The paper of Frees, Meyers and Cummings (2011a) also derived distribution theory to assess statistical differences between Gini indices. Although we do not review that theory here, we did perform these calculations for our data. It turns out that there is no statistically significant differences among the ten Gini indices that are based on the extended set of explanatory variables.

In summary, Table 3 suggests that there are important advantages to using extended sets of variables compared to the basic variables, regardless of the scoring techniques used.

### 3.3 Comparison of Scores Using the Gini Index

As demonstrated in the preceding section, if a base premium is available, then the Gini index can be used to decide whether an alternative insurance score is useful for detecting differences between loss and premium distributions. In instances where no base premium is available, the Gini index is also useful although care is required when interpreting this measure.

Table 4 summarizes the calculation of several Gini indices. Here, we allow the “base premium”  $P(\cdot)$  to be each of the fourteen competing scores plus the benchmark “ConsPrem,” a premium that is constant over policyholders. For each base premium, Table 4 shows the Gini index for each of the thirteen competing scores. Standard errors are reported in Table 5.

From the first row of Table 4, we see that all of the Gini indices are large, indicating that any of the fourteen scores considered here provide useful separation between losses and a constant premium. The next block (of four rows) consists of scores that use the basic set of explanatory variables, including SP\_FreqSev\_Basic. These four scores seem to perform similarly. For example, when compared to one another, the Gini indices are in the single digits. When any of the four are adopted as the base premium, double digit Gini indices are possible using the extended set of explanatory variables.

Using the Gini measure, the dependence ratio scores advanced by Frees, Meyers and Cummings (2010) seem to fare poorly. Almost every alternative score, except for the independence multi-peril frequency and severity models upon which they are based, allows for an ordering where there is a substantial separation between premium and loss distributions.

One approach for selecting a score based on the Gini index is a “mini-max” strategy. That is, select the score that provides the smallest of the maximal Gini indices, taken over competing scores. The strategy is intuitively appealing. If one were to specify a base premium, then the maximum Gini index corresponds to the largest separation between the loss and premium distribution when considering different orderings. For example, when the base premium is SP\_PurePrem\_Basic, the maximum Gini index is 12.8 which is achieved when IV\_FreqSevC is used to compute relativities to order distributions. For this criterion, Table 4 shows that IV\_FreqSevA is the best score. It has the smallest maximum Gini index at 7.2. We interpret this to mean that this score is the least vulnerable to alternative scores.

Table 5 shows that the standard errors of the Gini indices are relatively stable across different choices of scores and premiums. We will use this observation in Section 6 to propose some rules of thumb for sample size determination.

## 4 Using Covariances to Express the Gini Index

Equation (3) defines our Gini index in terms of an area associated with the ordered Lorenz curve. For the classic Lorenz curve and associated Gini index, there are several alternative (equivalent) definitions, cf., Yitzhaki (1998). These different definitions encourage alternative interpretations of the Gini index, hence widening the scope of potential applications. As with the classic Gini, we can also provide an alternative expression for our Gini index using covariance operators.

#### 4.1 The Gini Index in Terms of Covariances

We use the notation  $\widehat{\text{Cov}}(y, P)$  to denote the (empirical) covariance between losses  $y$  and premiums  $P$ . That is, define  $\widehat{\text{Cov}}(y, P) = n^{-1} (\sum_{i=1}^n y_i P_i - n\bar{y}\bar{P})$  (and recall that  $\bar{y} = \bar{P} = 1$ ). Then, after some pleasant algebra (see Appendix Section 9), we can express the Gini index as

$$\widehat{Gini} = 2\widehat{\text{Cov}}(y, \hat{F}_P(R)) - 2\widehat{\text{Cov}}(P, \hat{F}_R) - \frac{1}{n}\widehat{\text{Cov}}(y, P), \quad (5)$$

where  $\hat{F}_R = \text{rank}(R)/n$  is the distribution function of the rank of relativities. For large sample sizes  $n$ , the third term on the right-hand side of equation (5) is small and can be ignored.

With equation (4), we interpret a low relativity means that a policy is highly profitable and a good candidate to retain. Additional insights arise from equation (5). Other things being equal:

1. Under the relativity ordering, a large covariance between losses ( $y$ ) and the proportion of premiums retained ( $\hat{F}_P(R)$ ) implies a high Gini index.
2. A large negative covariance between premiums ( $P$ ) and relativities ( $\hat{F}_R$ ) implies a high Gini index. Stated differently, low relativities associated with high premiums implies a high Gini index. We retain policies with a low relativity. Other things being equal, it is more profitable to retain a policy with a high premium.

For many datasets, we have found that, using equation (5), we can approximate the weighted premium distribution  $\hat{F}_P(R)$  with the unweighted distribution of relativities  $\hat{F}_R$ . With this, we may define

$$\widehat{Gini}_{Approx} = \frac{2}{n}\widehat{\text{Cov}}((y - P), \text{rank}(R)). \quad (6)$$

Although this approximation to the Gini index provides little advantage computationally, it does give us another way to interpret the Gini index. We can think about  $P - y$  as the “profit” associated with a policy. Then, we may interpret the Gini index to be proportional to the negative covariance between profits and the rank of relativities. That is, if policies with low profits are associated with high relativities and high profits are associated with low relativities, then we have a profitable situation meaning that the Gini index is positive and large.

When “premiums”  $P$  are interpreted as exposures, it is more helpful to think in terms of pure premiums. An alternative approximation is

$$\widehat{Gini}_{Approx2} = \frac{2}{n}\widehat{\text{Cov}}(PP, \text{rank}(R)), \quad (7)$$

where the scaled pure premium  $PP$  is loss per premium ( $y/P$ ) rescaled (divided) by its average.

**Homeowners Example.** To understand the reliability of these approximations, we return to the homeowners example and show summaries in Table 6 and Figure 4. In Figure 4, we compare the two approximations to the Gini indices produced in Table 4. The left-hand panel shows that the approximation given in equation (6) to be very robust over a wide range of premiums and relativities. The right-hand panel shows that the approximation given in equation (7) to be less robust although still helpful for interpretation purposes.

In Table 6, we show the Gini indices and approximations for only those using SP\_FreqSev\_Basic as a base premium. We also decompose the equation (6) approximation into “loss” and “premium” sources. Here, the loss source is given by  $2\widehat{Cov}(y, rank(R))/n$  and similarly for premiums. These two sources can help interpret the magnitude of the Gini index.

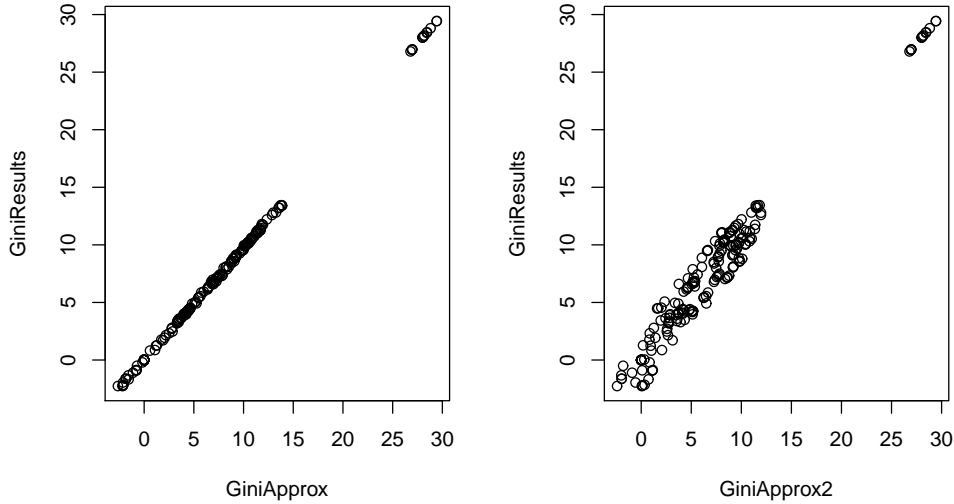


Figure 4: Gini Indices and Approximations. Observations in the upper-right hand corner correspond to using a constant premium as a base.

## 4.2 Some Special Cases

### 4.2.1 Simple Gini Index

In some applications, one can use an exposure measure, such as car years, instead of a premium. In others, it is helpful to think about the premium and exposure as constant over policies. In this

case, the relativity  $R(\mathbf{x}_i) = S(\mathbf{x}_i)$  is simply an insurance score. The direct comparison of losses to scores results in what we call a “simple Gini.”

When premiums are constant ( $P_i \equiv 1$ ), interpretations (and the algebra) are simpler. In this case, the relativity is the score ( $R_i = S_i$  in our notation). From equation (5), the Gini index reduces to

$$\widehat{Gini} = \frac{2}{n} \widehat{\text{Cov}}(y, \text{Rank}(S)). \quad (8)$$

Equation (8) is an exact relationship, no approximations are involved. It states that the simple Gini index is proportional to the covariance between losses and the rank of scores. Note that it is not a Pearson correlation between losses and scores, nor is it a Spearman correlation (the correlation between ranks of losses and ranks of scores). As discussed in Frees et al. (2011a), this statistic seems to have been first proposed by Durbin (1954) who proposed it as an instrumental variable estimator in an errors-in-variables regression problem. Durbin argued that using the rank of an explanatory variable may be helpful in explaining the behavior of  $y$  when values of the explanatory variable are mis-measured.

#### 4.2.2 Reverse Gini Index

We now reverse the roles of scores  $S$  and premiums  $P$  and call the resulting Gini index a “reverse Gini.” Returning to equation (6), an approximation for the reverse Gini is

$$\begin{aligned} \widehat{GiniR}_{Approx} &= \frac{2}{n} \widehat{\text{Cov}}((y - S), (n + 1 - \text{rank}(R))) \\ &= \frac{2}{n} \widehat{\text{Cov}}((S - y), \text{rank}(R)). \end{aligned} \quad (9)$$

Moreover, suppose that the score  $S$  is an unbiased estimator of the loss in the sense that  $E(y|\mathbf{x}) = S$ . Then,

$$\begin{aligned} \text{Cov}((S - y), \text{rank}(R)) &= E\{(S - y) \times \text{rank}(R)\} - E(S - y) \times E \text{rank}(R) \\ &= E\{(E(S - y)|\mathbf{x}) \times \text{rank}(R)\} - E(E(S - y)|\mathbf{x}) \times E \text{rank}(R) \\ &= 0, \end{aligned}$$

because  $E(y|\mathbf{x}) = S$ . This suggests that one can anticipate the reverse Gini,  $\widehat{GiniR}_{Approx}$ , to be zero when the model is well-specified. We use the reverse Gini as another statistic to measure model fit.

### 4.2.3 Gini Index for Refined Scores

For another case, suppose that the score  $S$  is a “more refined” version of a premium  $P$ . For example,  $S$  may reflect information in a new rating variable (such as a credit score) or more precise geographic information. Specifically, let  $S = P \exp(\mathbf{z}'\boldsymbol{\beta})$ , where  $\mathbf{z}$  is a vector of new variables not contained in the premium base  $P$ . It is helpful to think about some specific examples.

**Continuous Variable.** Suppose that we consider on a single continuous variable (e.g., credit score). Then, the rank of the relative premium can be expressed as

$$\text{rank}(R) = \text{rank}\left(\frac{S}{P}\right) = \text{rank}(\exp(z\beta)) = \text{rank}(z),$$

assuming that  $\beta$  is positive. Then, from equation (6), we may interpret the Gini index to be approximately

$$\widehat{Gini}_{Approx} = \frac{2}{n} \widehat{\text{Cov}}((y - P), \text{rank}(z)),$$

that is, the Gini is approximately the covariance between the policy “profit”  $P - y$  and the rank of the new variable  $z$ , rescaled by the constants.

**Categorical Variable.** Suppose that we consider on a single discrete variable  $z$  with three possible outcomes 1, 2, and 3 (e.g., urban, suburb and rural). Suppose we use  $\mathbf{z}'\boldsymbol{\beta} = \beta_1 \mathbf{I}(z = 1) + \beta_2 \mathbf{I}(z = 2) + \beta_3 \mathbf{I}(z = 3)$ . Here, recall that  $\mathbf{I}(\cdot)$  is the indicator function. Without loss of generality, assume that  $\beta_1 < \beta_2 < \beta_3$  (otherwise, simply re-order  $z$ ). Then,  $\text{rank}(R) = \text{rank}(\exp(z\beta)) = \text{rank}(z) = z$ , and

$$\widehat{Gini}_{Approx} = \frac{2}{n} \widehat{\text{Cov}}((y - P), z).$$

Each level of  $z$  represents a “segment” of the market that is implemented in the new score  $S$  but not in the original premium base  $P$ . The Gini index measures the relationship between the policy “profit”  $P - y$  and the market segments in  $z$ .

For both examples, we can think about the Gini index as summarizing the linear relationship between policy profit  $P - y$  and the rank of the refinement variable,  $\text{rank}(z)$ . This suggests additional analyses, such as a plot of  $P - y$  versus  $\text{rank}(z)$ , in order to understand potential nonlinear relationships.

## 5 Model Selection

In this section, we investigate the role of the Gini index as a statistic to aid in selecting a model through a simulation study.



## 5.1 Simulation Study Design

The study is designed to replicate many of the data features that we encountered when analyzing the homeowners data described in Section 3.

For each scenario, we generated  $n$  in-sample policyholder observations, estimated model parameters and then calculated scores for each of  $n$  out-of-sample policyholders. In our simulation, we let  $n$  equal 500,000.

**Simulated Distributions.** For each policyholder, we assumed knowledge of two characteristics where each  $x_j$  was generated from a chi-square distribution with 20 degrees of freedom, rescaled to have a zero mean and variance 1/10. With these choices, the score distributions (score calculations are described below) exhibited a right-skewed distribution comparable to the premium distribution portrayed in Figure 1. The regression function was generated using a logarithmic link function, that is,  $m(\mathbf{x}) = \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$ . Using the regression function as the location parameter, we generated the loss  $y$  using the Tweedie distribution. Here, parameters of the Tweedie distribution were set so that the simulated distribution was comparable to the distribution of our homeowners data in Section 3. We did this for  $n$  in-sample and  $n$  out-sample observations, respectively.

**Score Calculation.** Using the in-sample data, we estimated parameters for each of eight scores:

- $S_1(\mathbf{x}) = \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2)$ , the true regression function
- $S_2(\mathbf{x}) = \exp(\beta_0 + \beta_1 x_1)$ , based on only  $x_1$
- $S_3(\mathbf{x}) = \exp(\beta_0 + \beta_2 x_2)$ , based on only  $x_2$
- $S_4(\mathbf{x}) = \exp(\beta_0 + \beta_1 \frac{1}{x_1})$ , based on an (incorrect) reciprocal transform of  $x_1$
- $S_5(\mathbf{x}) = \exp(\beta_0 + \beta_2 \frac{1}{x_2})$ , based on an (incorrect) reciprocal transform of  $x_2$
- $S_6(\mathbf{x}) = \exp(\beta_0 + \beta_1 x_1 + \beta_2 \frac{1}{x_2})$ , based on  $x_1$  and an (incorrect) reciprocal transform of  $x_2$
- $S_7(\mathbf{x}) = \exp(\beta_0 + \beta_1 \frac{1}{x_1} + \beta_2 x_2)$ , based on  $x_2$  and an (incorrect) reciprocal transform of  $x_1$
- $S_8(\mathbf{x}) = \exp(\beta_0 + \beta_1 \frac{1}{x_1} + \beta_2 \frac{1}{x_2})$ , based on (incorrect) reciprocal transforms of  $x_1$  and  $x_2$

Then, with the parameter estimates of the score coefficients from the in-sample data, we used the out-of-sample characteristics to generate scores. We then compared scores to the actual out-of-sample losses. We report results based on 100 simulations. With these sample and simulation sizes, it turns out that the simulation standard errors for all Gini indices are less than 0.2%.

Using the regression function in score 1, we generated three scenarios by varying the regression coefficients. In the first scenario, we let  $\beta_1 = \beta_2 = 0.25$  so that both explanatory variables contribute equally to the regression function. In the second, we let  $\beta_1 = 0.25$  and  $\beta_2 = 0.05$  so that the second explanatory variable contributes little to the regression function. In the third scenario, we let  $\beta_1 = 0.05$  and  $\beta_2 = 0.25$  so that the first explanatory variable contributes little to the regression function.

For each scenario, we assume that the analyst is considering one of eight scores (with the first being the correct choice, unknown to the analyst). To measure the discrepancy between the chosen score and the true regression function, we present the Spearman correlation that we label as the “True Correlation.” Note that even when the analyst chooses the correct score, the Spearman correlation is less than one due to the (in-sample) estimation error in the regression coefficients in the score. If the true correlation were available, then the analyst would simply choose the model with the largest true correlation. However, this is unavailable, and so the analyst must use available statistics, including the “Simple Gini” and the ratio Gini indices presented in Table 7. An analyst could select a model by searching for the largest simple Gini index. Alternatively, one could use a mini-max strategy discussed in Section 3.3.

## 5.2 Simulation Study Results

Table 7 summarizes the results of the simulation study. Note that through our choice of scenarios we may observe outcomes over a broad range of models selected, ranging from a near perfect selection (where the correlation is 99.29%) to a very poor selection (where the correlation is only 5.35%).

The simple Gini index seems to be a desirable proxy for the true correlation. Over different scenarios and different premiums, as the simple Gini index increases, so does the true correlation. This is intuitive plausible in that the simple Gini index may be interpreted as proportional to the covariance between the insurance loss and rank of the score and whereas the “true correlation” is the correlation between the rank of the regression function and the rank of scores (a Spearman correlation). If one converts the simple Gini to a correlation it turns out to be much smaller than the true correlation. This is simply because of the noise in the loss random variable as a estimate of its expectation, the regression function.

Choosing the smallest of the maximum Gini ratios is also a viable model selection strategy. Table 7 shows a strong inverse relation between the “maximum” column and the true correlation column.

The simulation also allows us to document the “reverse Gini effect.” To see this, consider the first scenario and suppose that the analyst initially chooses Score 2 as the base premium. For Score 3 as an alternative, the resulting ratio Gini index is 4.97% suggesting that this score is preferred. However, if the analysts using Score 3 as the base and Score 2 as the alternative, then Table 7 shows that the resulting Gini index is 5.12, suggesting that Score 2 is preferred. This is the reverse Gini effect, where the Gini analysis provides seemingly contradictory advice.

However, because we generated the scores and the model, we know that neither Score 2 nor Score 3 represents the true outcome. Thus, the decision-making process with Gini indices suffers from the same drawback as with statistical hypothesis testing. Model A can be rejected in favor of Model B and vice-versa if neither model is true. Table 7 shows that the reverse Gini effect is not present in the second and third scenarios when one variable dominates the other.

## 6 Sample Size Determination

How large a sample size is required for a reliable *Gini* statistic? In this section, we show how to use results from the theoretical properties, plus some basic knowledge of the loss distribution, to provide rules of thumb that can be used to select an appropriate sample size. This procedure could be used, for example, to determine the size of a block of business when examining alternative premium structures on a trial, or “pilot,” basis.

In earlier work, we showed that the distribution of the *Gini* statistic is approximately normal for large samples, see Theorem A2 of Appendix Section 8. The form of the large sample variance,  $\Sigma_{Gini}/n$ , given in Theorem A2, is complicated. However, Appendix Section 10 shows that using the assumption of independent relativities results in a much simpler expression, given as

$$\Sigma_{Gini} = \frac{\text{Var}(y - P)}{3}. \quad (10)$$

This result is similar to one for the Pearson correlation, another measure of association, where the form of the variance simplifies under the independence assumption.

To illustrate, for our sample described in Section 3, we have  $n = 359,454$  observations. After rescaling so that premiums and losses are mean one, we have the standard deviation of losses and premiums are, respectively,  $s_y = 14.79591$  and  $s_P = 0.70558$ . The covariance between losses and premium is  $Cov_{yP} = 0.48538$ . From this and equation (10), an approximate standard error for the

Gini index is

$$se(\widehat{Gini})_{Approx} = \sqrt{\frac{14.79591^2 + 0.70558^2 - 2 \times 0.48538}{3 \times 359454}} = 0.0142.$$

This result is close to the standard errors presented in Table 5 (calculated using the complicated, yet more precise, expression for  $\Sigma_{Gini}$  given in Appendix Section 8).

Because this approximation is valid over a range of relativities and premiums (with different dependency structures), we conclude that this approximation is helpful for determining the size of a sample to be collected for test studies. Figure 5 shows the effect of alternative sample sizes on the approximate Gini standard error for our homeowners data.

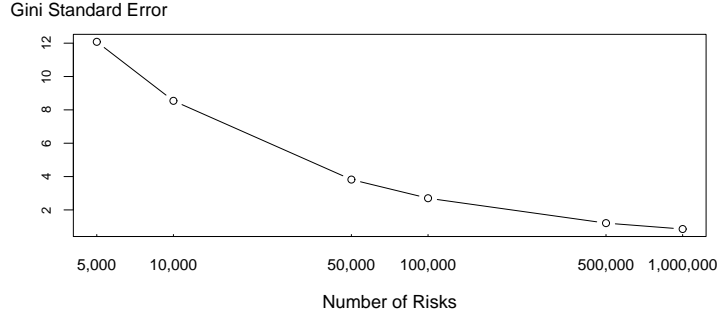


Figure 5: Effect of Sample Size on Gini Approximate Standard Errors.

## 7 Summary and Concluding Remarks

The Gini index is a measure of association between losses and premiums - one that has important economic content in insurance scoring applications. For a given ordering of risks, the Gini index summarizes the difference between the premium and loss distributions. An excess of premiums over losses can be interpreted to be an insurer's profit. This observation leads an insurer to seek an ordering that produces to a large Gini index. Thus, the Gini index and associated ordered Lorenz curve are useful for identifying profitable blocks of insurance business.

Unlike classical measures of association, the Gini index assumes that a premium base  $P$  is currently in place and seeks to assess vulnerabilities of this structure. This approach is more akin to hypothesis testing (when compared to goodness of fit) where one identifies a “null hypothesis” as the current state of the world and uses decision-making criteria/statistics to compare this to an

“alternative hypothesis.” The purpose of this paper is not to say that either hypothesis testing or goodness of fit approaches are always good or bad; rather, both have their place in statistical inference. The purpose of the paper is provide another measure that can be used to augment the analyst’s toolkit; we argue that this new measure provides insights that are not available from classical measures of association.

To summarize, we anticipate the Gini index being of use in at least the following three situations:

1. A premium structure  $P$  is in place and we wish to assess the usefulness of a generic alternative score  $S$ . This is the basic scenario in which we introduced the ordered Lorenz curve and the Gini index to combat adverse selection. We also discussed the Gini index as a measure of profit in Section 2. To illustrate, we demonstrated its usefulness in the homeowners example in Section 3.2; analysts can supplement this analysis by looking to the reverse Gini as an additional measure of model fit.
2. A premium structure  $P$  is in place and the alternative score is a refined version of the premium. Although this is the same as the basic scenario defined above, additional interpretations are available for the relativity that are potentially helpful in model diagnostics.
3. No premium structure is in place - a number of alternative scoring methods are being considered. In this case, at least two strategies are available. One is to use the “minimax” strategy put forth in Section 3.3, where one chooses a score that is least vulnerable to competition from other scores. The other strategy is to use the “simple” Gini index that has no base premium as a reference.

To assess the reliability of the Gini index, Section 6 describes principles for sample size determination. On the one hand, sample sizes required for reliable applications such as in personal lines homeowners insurance are quite large, in some cases ranging into the hundreds of thousands of observations. On the other hand, with large sample sizes available, we can enjoy reliable inferences for complex distributions that are mixtures of a large mass at zero and a right-skewed, thick-tailed positive distribution. Given the availability of large datasets in today’s world, we view these sample size requirements as feasible in some important areas of applications.

As a concrete illustration, this paper has focussed on insurance as a financial risk. Further, Frees et al. (2011a) showed that traditional credit scoring arises as a special case of our formulation. Moreover, additional potential applications in credit scoring are easy to imagine. For example, one could let  $y$  represent the *amount* of credit default (not just the event) and allow the amount charged

for the loan to depend on an applicant's creditworthiness. Results of this paper apply directly to this situation.

## References

- Durbin, J. (1954). Errors in variables. *Review of International Statistical Institute* 22, 23-32.
- Egghe, L. (2005). Continuous, weighted Lorenz theory and applications to the study of fractional relative impact factors. *Information Processing and Management* 41, 1330-1359.
- Frees, Edward W., Glenn Meyers and A. David Cummings (2010). Dependent multi-peril ratemaking models. *Astin Bulletin* 40(2), 699-726.
- Frees, Edward W., Glenn Meyers and A. David Cummings (2011a). Summarizing insurance scores using a Gini index. To appear, *Journal of the American Statistical Association*. Available at <http://research3.bus.wisc.edu/jfrees>.
- Frees, Edward W., Glenn Meyers and A. David Cummings (2011b). Predictive modeling of multi-peril homeowners insurance. Working paper. Available at <http://research3.bus.wisc.edu/jfrees>.
- Gastwirth, Joseph L. (1971). A general definition of the Lorenz curve. *Econometrica* 39 (6), 1037-1039.
- Gastwirth, Joseph L. (1972). The estimation of the Lorenz curve and the Gini index. *Review of Economics and Statistics* 54 (3), 306-316.
- Gastwirth, Joseph L. (1975). Statistical measures of earnings differentials. *The American Statistician* 29 (1), 32-35.
- Hastie, Trevor, Robert Tibshirani and Jerome Friedman (2001). *The Elements of Statistical Learning: Data Mining, Inference and Prediction*. Springer, New York.
- Lambert, Peter J. and André Decoster (2005). The Gini coefficient reveals more. *Metron* 63 (3), 373-400.
- Lorenz, Max O. (1905). Methods of measuring the concentration of wealth. *Publications of the American Statistical Association* 9 (70), 209-219.
- Meyers, Glenn and A. David Cummings (2009). "Goodness of Fit" vs. "Goodness of Lift". *The Actuarial Review: Newsletter of the Casualty Actuarial Society*, August, p. 16. Available at [http://www.casact.org/newsletter/pdfUpload/ar/AR\\_Aug2009\\_1.pdf](http://www.casact.org/newsletter/pdfUpload/ar/AR_Aug2009_1.pdf).
- Nicodemus, Kristin K. and James D. Malley (2009). Predictor correlation impacts machine learning algorithms: implications for genomic studies. *Bioinformatics* 25 (15), 1884-1890.
- Sandri, Marco and Paola and Zuccolotto (2008). A bias correction algorithm for the Gini variable importance measure in classification trees. *Journal of Computational and Graphical Statistics* 17 (3), 611-628.
- Sen, Amartya and James E. Foster (1998). *On Economic Inequality*. Oxford University Press, Delhi.
- UNU-WIDER World Income Inequality Database (2008). Version 2.0c, May 2008. Sponsored by the World Institute for Development Economics Research of the United Nations University, Helsinki, Finland. Available at: [http://www.wider.unu.edu/research/Database/en\\_GB/database/](http://www.wider.unu.edu/research/Database/en_GB/database/).
- Yitzhaki, Shlomo (1996). On using linear regressions in welfare economics. *Journal of Business and Economic Statistics* 14 (4), 478-486.
- Yitzhaki, Shlomo (1998). More than a dozen ways of spelling Gini. *Research on Economic Inequality* 8, 13-30.

## 8 Appendix A - Properties of the Gini Index

The (population) *Gini* index is

$$Gini = 2 \int_0^{\infty} \{F_P(s) - F_L(s)\} dF_P(s). \quad (11)$$

Here,  $F_P$  and  $F_L$  are weighted distributions functions which are the population versions of the empirical distributions given in equations (1) and (2), respectively.

We summarize the consistency and asymptotic normality of the empirical Gini index  $\widehat{Gini}$  in the following two results.

**Theorem A1.** Under mild regularity conditions, the Gini statistic  $\widehat{Gini}$  is a consistent estimator of the Gini index. That is,  $\widehat{Gini} \rightarrow Gini$ , as  $n \rightarrow \infty$ , with probability one.

For asymptotic normality, we use the projection

$$h_1(\mathbf{x}, y) = \frac{1}{2} (\mu_y P(\mathbf{x}) F_L(R) + y \mu_P [1 - F_P(R)]). \quad (12)$$

Further, use the notation  $\Sigma_h = \text{Var } h_1(\mathbf{x}, y)$ ,  $\Sigma_y = \text{Var } y$ ,  $\Sigma_P = \text{Var } P(\mathbf{x})$ ,  $\Sigma_{hy} = \text{Cov}(h_1(\mathbf{x}, y), y)$ ,  $\Sigma_{yP} = \text{Cov}(y, P(\mathbf{x}))$ , and  $\Sigma_{hP} = \text{Cov}(h_1(\mathbf{x}, y), P(\mathbf{x}))$ . With these terms, we can establish:

**Theorem A2.** Under mild regularity conditions, the Gini statistic  $\widehat{Gini}$  has an asymptotic normal distribution. Specifically,  $\sqrt{n} (\widehat{Gini} - Gini) \rightarrow_D N(0, \Sigma_{Gini})$ , where

$$\Sigma_{Gini} = \frac{4}{\mu_y^2 \mu_P^2} \left( 4\Sigma_h + \frac{\mu_h^2}{\mu_y^2} \Sigma_y + \frac{\mu_h^2}{\mu_P^2} \Sigma_P - \frac{4\mu_h}{\mu_y} \Sigma_{hy} - \frac{4\mu_h}{\mu_P} \Sigma_{hP} + \frac{2\mu_h^2}{\mu_y \mu_P} \Sigma_{yP} \right), \quad (13)$$

with  $\mu_h = \mu_y \mu_P (1 - Gini)/2$ .

To estimate the asymptotic variance, Table 1 provides moment-based estimators.

Table 1: Moment-Based Estimators for the Asymptotic Variance	
$\hat{h}_1(\mathbf{x}, y) = \frac{1}{2} (P(\mathbf{x}) \hat{F}_L(R) + y[1 - \hat{F}_P(R)])$	$\hat{\Sigma}_{hP} = n^{-1} \sum_{i=1}^n \hat{h}_1(\mathbf{x}_i, y_i) P(\mathbf{x}_i) - \bar{h}_1$
$\bar{h}_1 = n^{-1} \sum_{i=1}^n \hat{h}_1(\mathbf{x}_i, y_i)$	$\hat{\Sigma}_y = n^{-1} \sum_{i=1}^n y_i^2 - 1$
$\hat{\Sigma}_h = n^{-1} \sum_{i=1}^n \hat{h}_1(\mathbf{x}_i, y_i)^2 - \bar{h}_1^2$	$\hat{\Sigma}_P = n^{-1} \sum_{i=1}^n P(\mathbf{x}_i)^2 - 1$
$\hat{\Sigma}_{hy} = n^{-1} \sum_{i=1}^n \hat{h}_1(\mathbf{x}_i, y_i) y_i - \bar{h}_1$	$\hat{\Sigma}_{yP} = n^{-1} \sum_{i=1}^n y_i P(\mathbf{x}_i) - 1$
<i>Note:</i> These estimators are based on rescaling so that $\bar{y} = \bar{P} = 1$ .	

**Theorem A3.** Under mild regularity conditions, a consistent estimator of  $\Sigma_{Gini}$  is

$$\hat{\Sigma}_{Gini} = 4 \left( 4\hat{\Sigma}_h + \bar{h}_1^2 \hat{\Sigma}_y + \bar{h}_1^2 \hat{\Sigma}_P - 4\bar{h}_1 \hat{\Sigma}_{hy} - 4\bar{h}_1 \hat{\Sigma}_{hP} + 2\bar{h}_1^2 \hat{\Sigma}_{yP} \right). \quad (14)$$

The proofs of these results are Frees, Meyers and Cummings (2011a).

## 9 Appendix B - Proof of Equation (5)

We establish equation (5) assuming  $\bar{y} = \bar{P} = 1$ . First, from equation (1), we have that  $a_j - a_{j-1} = \frac{P_j}{n}$  and, from equation (2), we have that  $b_j + b_{j-1} = 2\hat{F}_L(R_j) - \frac{y_j}{n}$ . Thus, with equation (3), we have

$$\begin{aligned}\widehat{Gini} &= 1 - \sum_{j=1}^n (a_j - a_{j-1})(b_j + b_{j-1}) \\ &= 1 - \frac{1}{n} \sum_{j=1}^n P_j \left( 2\hat{F}_L(R_j) - \frac{y_j}{n} \right).\end{aligned}\tag{15}$$

Now, with equation (2) and a change of summations, we can write

$$\begin{aligned}\sum_{j=1}^n P_j \hat{F}_L(R_j) &= \frac{1}{n} \sum_{j=1}^n \sum_{i=1}^n y_i P_j \mathbf{I}(R_i \leq R_j) \\ &= \frac{1}{n} \sum_{i=1}^n y_i \left\{ \sum_{j=1}^n P_j \mathbf{I}(R_i \leq R_j) \right\} \\ &= \frac{1}{n} \sum_{i=1}^n y_i \left\{ \sum_{j=1}^n P_j (1 - \mathbf{I}(R_j \leq R_i) + \mathbf{I}(R_i = R_j)) \right\} \\ &= \frac{1}{n} \sum_{i=1}^n y_i \left\{ n - n\hat{F}_P(R_i) + P_i \right\}.\end{aligned}$$

Putting this into equation (15), we have

$$\begin{aligned}\widehat{Gini} &= 1 - \frac{2}{n^2} \sum_{j=1}^n y_i \left\{ n - n\hat{F}_P(R_i) + P_i \right\} + \frac{1}{n^2} \sum_{i=1}^n y_i P_i \\ &= \frac{2}{n} \sum_{j=1}^n y_i \hat{F}_P(R_i) - 1 - \frac{1}{n^2} \sum_{i=1}^n y_i P_i \\ &= \frac{2}{n} \left\{ n\widehat{\text{Cov}}(y, \hat{F}_P(R)) + n\overline{\hat{F}_P(R)} \right\} - 1 - \frac{1}{n^2} \left\{ n\widehat{\text{Cov}}(y, P) + n \right\} \\ &= 2\widehat{\text{Cov}}(y, \hat{F}_P(R)) + 2\overline{\hat{F}_P(R)} - \frac{n+1}{n} - \frac{1}{n} \widehat{\text{Cov}}(y, P).\end{aligned}\tag{16}$$

We now use  $\widehat{\text{Cov}}(P, R) = n^{-1}(\sum_{i=1}^n P_i \times i - n\frac{n+1}{2})$  (recall that premiums are sorted by relativities so that the rank of the  $i$ th relativity is  $i$ ). To calculate the average weighted premium distribution,



we have

$$\begin{aligned}
\overline{\hat{F}_P(R)} &= \frac{1}{n} \sum_{i=1}^n \hat{F}_P(R_i) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n P_j \mathbf{I}(R_j \leq R_i) \\
&= \frac{1}{n^2} \sum_{j=1}^n P_j (n - j + 1) \\
&= \frac{n+1}{n} - \frac{1}{n^2} \left\{ n \widehat{\text{Cov}}(P, R) + n \frac{n+1}{2} \right\} \\
&= \frac{n+1}{n} - \frac{n+1}{2n} - \frac{1}{n} \widehat{\text{Cov}}(P, R) = \frac{n+1}{2n} - \widehat{\text{Cov}}(P, F_R).
\end{aligned}$$

Putting this into equation (16) yields

$$\begin{aligned}
\widehat{Gini} &= 2\widehat{\text{Cov}}(y, \hat{F}_P(R)) + 2\left(\frac{n+1}{2n} - \widehat{\text{Cov}}(P, F_R)\right) - \frac{n+1}{n} - \frac{1}{n} \widehat{\text{Cov}}(y, P) \\
&= 2\widehat{\text{Cov}}(y, \hat{F}_P(R)) - 2\widehat{\text{Cov}}(P, F_R) - \frac{1}{n} \widehat{\text{Cov}}(y, P),
\end{aligned}$$

which is equation (5).  $\square$

## 10 Appendix C - Sample Size Calculations

The following proposition is a corollary of Theorem A2 of Frees, Meyers and Cummings (2011a).

**Proposition.** Assume that  $R$  is independent of  $(y, P)$  and the conditions of Theorem A2 hold. Then, we have  $\sqrt{n} \widehat{Gini} \rightarrow_D N(0, \Sigma_{Gini})$ , where

$$\Sigma_{Gini} = \frac{\Sigma_y + \Sigma_P - 2\Sigma_{yP}}{3} = \frac{\text{Var}(y - P)}{3}.$$

For simplicity, we establish this proposition assuming by losses and premiums have been rescaled by dividing by their respective averages. Through this rescaling, we have that the mean loss is  $E y = 1$  and the mean premium is  $E P = 1$ .

*Proof.* Under the assumption that  $R$  is independent of  $(y, P)$ , we first note that

$$\begin{aligned}
F_P(s) &= \frac{E[P(\mathbf{x})\mathbf{I}(R \leq s)]}{E P(\mathbf{x})} = \frac{E[P(\mathbf{x})] \Pr(R \leq s)}{E P(\mathbf{x})} \\
&= \Pr(R \leq s) = F_R(s),
\end{aligned}$$

and similarly,  $F_L(s) = F_R(s)$ . Thus, using equation (12), we may write the projection as

$$h_1(\mathbf{x}, y) = \frac{1}{2} (PF_R + y[1 - F_R]).$$

Recall, for continuous relativities  $R$ , that  $F_R$  has a uniform distribution and so  $E F_R = 1/2$ ,  $\text{Var } F_R = 1/12$ , and  $E F_R^2 = 1/12 + (1/2)^2 = 1/3$ . Now, assuming  $R$  is independent of  $(y, P)$ , this has mean

$$\mu_h = E h_1(\mathbf{x}, y) = \frac{1}{2} \{ (1)(1/2) + (1)[1 - (1/2)] \} = \frac{1}{2},$$

and so  $Gini = 0$ . Further,

$$\begin{aligned} \Sigma_h &= \text{Var } h_1(\mathbf{x}, y) = E h_1^2 - \left( \frac{1}{2} \right)^2 \\ &= \frac{1}{4} E \{ P^2 F_R^2 + y^2 (1 - F_R)^2 + 2yP(F_R - F_R^2) \} - \frac{1}{4} \\ &= \frac{1}{4} \left( (\Sigma_P + 1) \left( \frac{1}{3} \right) + (\Sigma_y + 1) \left( \frac{1}{3} \right) + 2(\Sigma_{yP} + 1) \left( \frac{1}{2} - \frac{1}{3} \right) \right) - \frac{1}{4} \\ &= \frac{\Sigma_y + \Sigma_P + \Sigma_{yP}}{12}. \end{aligned}$$

Similarly,

$$\begin{aligned} \Sigma_{hy} &= \text{Cov } h_1(\mathbf{x}, y) = E y h_1 - \left( (1) \frac{1}{2} \right) \\ &= \frac{1}{2} E \{ y P F_R + y^2 (1 - F_R) \} - \frac{1}{2} \\ &= \frac{1}{2} E \left\{ (\Sigma_{yP} + 1) \frac{1}{2} + (\Sigma_y + 1) \frac{1}{2} \right\} - \frac{1}{2} \\ &= \frac{\Sigma_y + \Sigma_{yP}}{4}. \end{aligned}$$

By symmetry, we have  $\Sigma_{hP} = (\Sigma_P + \Sigma_{yP})/4$ .

Now, using equation (13),  $\mu_y = 1$ ,  $\mu_P = 1$ , and  $\mu_h = 1/2$ , we have

$$\begin{aligned} \Sigma_{Gini} &= 4 \{ 4\Sigma_h + \mu_h^2 \Sigma_y + \mu_h^2 \Sigma_P - 4\mu_h \Sigma_{hy} - 4\mu_h \Sigma_{hP} + 2\mu_h^2 \Sigma_{yP} \} \\ &= 16\Sigma_h + \Sigma_y + \Sigma_P - 8\Sigma_{hy} - 8\Sigma_{hP} + 2\Sigma_{yP} \\ &= 16 \left( \frac{\Sigma_y + \Sigma_P + \Sigma_{yP}}{12} \right) + \Sigma_y + \Sigma_P - 8 \left( \frac{\Sigma_y + \Sigma_{yP}}{4} \right) - 8 \left( \frac{\Sigma_P + \Sigma_{yP}}{4} \right) + 2\Sigma_{yP} \\ &= \frac{\Sigma_y + \Sigma_P - 2\Sigma_{yP}}{3}, \end{aligned}$$

as required.  $\square$

Table 2: Summary Statistics of Fourteen Scores and Total Claims

Score	Percentiles									
	Mean	Mini- mum	1st	5th	25th	50th	75th	95th	99th	Maxi- mum
SP_FreqSev_Basic	291.10	20.48	85.00	120.25	182.74	240.37	334.62	618.37	1,025.88	8,856.79
SP_PurePrem_Basic	289.91	33.01	89.48	127.80	189.87	246.44	329.79	586.33	1,050.15	5,467.41
IND_PurePrem_Basic	290.91	37.49	92.08	124.04	182.68	240.30	328.87	612.47	1,087.06	13,577.91
IV_PurePrem_Basic	293.55	36.61	93.91	128.21	187.57	241.29	327.75	616.05	1,122.84	15,472.82
SP_FreqSev	287.79	8.78	71.55	105.39	171.55	237.95	339.40	631.98	1,039.19	6,864.46
SP_PurePrem	290.00	10.23	72.17	107.90	175.83	242.17	338.64	616.64	1,113.73	7,993.52
IND_FreqSev	294.93	33.05	97.14	126.61	185.07	244.99	333.68	606.03	1,106.17	22,402.49
IND_PurePrem	292.18	28.04	86.53	119.74	181.22	240.52	326.60	592.07	1,078.25	49,912.59
IV_PurePrem	294.06	12.42	78.41	113.14	178.62	240.38	330.21	614.22	1,095.70	107,158.09
IV_FreqSevA	290.91	23.99	88.70	121.70	182.29	241.42	327.81	606.23	1,096.86	18,102.93
IV_FreqSevB	295.32	28.52	94.58	124.77	184.29	245.26	335.38	606.63	1,100.61	24,394.06
IV_FreqSevC	291.17	20.88	84.78	118.21	180.63	241.57	329.92	608.28	1,098.40	20,046.03
DepRatio1	301.12	33.38	98.80	128.95	188.73	249.97	340.64	619.79	1,129.96	23,255.94
DepRatio36	302.39	33.48	99.27	129.65	189.87	251.41	342.30	620.38	1,132.36	23,092.35
TotClaims	332.89	0.00	0.00	0.00	0.00	0.00	0.00	660.00	5,916.33	350,000.00

Legend:

Score	Interpretation
<i>Scores using the basic set of explanatory variables</i>	
SP_FreqSev_Basic	Single-peril, frequency and severity model
SP_PurePrem_Basic	Single-peril, pure premium model
IND_PurePrem_Basic	Multi-peril independence, pure premium model
IV_PurePrem_Basic	Instrumental variable multi-peril pure premium model
<i>Scores using the extended set of explanatory variables</i>	
SP_FreqSev	Single-peril, frequency and severity model
SP_PurePrem	Single-peril, pure premium model
IND_FreqSev	Multi-peril frequency and severity model assuming independence among perils
IND_PurePrem	Multi-peril pure premium model assuming independence among perils
IV_PurePrem	Instrumental variable multi-peril pure premium model.
<i>Instrumental variable multi-peril frequency and severity models, using the extended set of explanatory variables</i>	
IV_FreqSevA	Uses instruments in frequency model
IV_FreqSevB	Uses instruments in severity model
IV_FreqSevC	Uses instruments in frequency and severity models
<i>Dependence ratio multi-peril frequency and severity models, using the extended set of explanatory variables</i>	
DepRatio1	Uses a single parameter for frequency dependencies
DepRatio36	Uses 36 parameters for frequency dependencies

Table 3: Gini Indices and Standard Errors

Alternative Score	Gini	Standard Error	Alternative Score	Gini	Standard Error
SP_PurePrem_Basic	4.89	1.43	IV_FreqSevA	12.59	1.39
IND_PurePrem_Basic	4.01	1.46	IV_FreqSevB	10.61	1.44
IV_PurePrem_Basic	4.33	1.46	IV_FreqSevC	12.80	1.38
SP_FreqSev	11.15	1.42	DepRatio1	10.09	1.45
SP_PurePrem	9.97	1.42	DepRatio36	10.06	1.46
IND_FreqSev	10.03	1.45			
IND_PurePrem	10.96	1.45			
IV_PurePrem	11.29	1.45			

Note: Base Premium is SP\_FreqSev\_Basic.

Table 4: Gini Indices for Fourteen Scores

Base Premium	Basic Explanatory Variables						Extended Explanatory Variables												Mazimum						
	Single Peril			IND_			IV_			Single Peril			IND_			IV_				DepRatio					
	Freq	Pure	Prem	Freq	Pure	Prem	Freq	Pure	Prem	Freq	Pure	Prem	Freq	Pure	Prem	Freq	Pure	Prem			A	B	C		
ConstPrem	27.0	27.0	26.8	26.8	26.9	26.9	28.8	28.1	28.0	28.0	28.5	28.4	28.4	29.4	28.2	29.4	28.0	28.0	29.4	28.0	28.0	29.4	28.0	28.0	29.4
SP_FreqSev_Basic	0.0	4.9	4.0	4.3	4.3	4.3	11.1	10.0	10.0	10.0	11.0	11.3	11.3	12.6	10.6	12.8	10.1	10.1	12.6	10.6	12.8	10.1	10.1	12.8	12.8
SP_PurePrem_Basic	4.3	0.0	2.7	3.3	3.3	3.3	11.2	8.0	8.0	8.0	9.9	9.9	9.9	11.4	8.8	11.7	8.1	8.1	11.4	8.8	11.7	8.1	8.1	11.7	11.7
IND_PurePrem_Basic	8.1	7.1	0.0	3.4	3.4	3.4	13.4	11.1	10.0	10.0	11.8	12.2	12.2	13.4	10.6	13.4	10.0	10.0	13.4	10.6	13.4	10.0	10.0	13.4	13.4
IV_PurePrem_Basic	7.9	6.6	3.6	0.0	0.0	0.0	12.8	10.7	10.2	10.2	11.5	11.7	11.7	13.2	10.6	13.2	10.2	10.2	13.2	10.6	13.2	10.2	10.2	13.2	13.2
SP_FreqSev	1.7	4.0	4.2	4.9	4.9	4.9	0.0	4.4	7.2	9.3	9.3	9.5	9.5	9.2	7.3	9.1	7.2	7.2	9.2	7.3	9.1	7.2	7.2	9.5	9.5
SP_PurePrem	7.0	5.9	6.8	7.0	7.0	7.0	9.1	0.0	8.6	9.7	9.7	9.5	9.5	10.3	8.8	10.5	8.6	8.6	10.3	8.8	10.5	8.6	8.6	10.5	10.5
IND_FreqSev	6.8	6.3	3.2	4.1	4.1	4.1	11.3	9.0	0.0	9.6	11.1	11.1	11.1	10.5	4.4	10.3	2.5	2.3	10.5	4.4	10.3	2.5	2.3	11.3	11.3
IND_PurePrem	6.1	5.0	1.2	2.2	2.2	2.2	8.6	6.8	4.2	0.0	3.7	3.7	3.7	7.4	4.2	7.3	4.3	4.2	7.4	4.2	7.3	4.3	4.2	8.6	8.6
IV_PurePrem	6.7	6.4	3.6	3.5	3.5	3.5	8.4	6.6	5.4	4.1	0.0	0.0	0.0	7.2	5.5	7.5	5.4	5.4	7.2	5.5	7.5	5.4	5.4	8.4	8.4
IV_FreqSevA	2.8	1.3	-1.1	-0.9	-0.9	-0.9	7.2	4.0	-2.3	4.5	5.1	5.1	5.1	0.0	-2.2	1.9	-2.2	-2.2	0.0	-2.2	1.9	-2.2	-2.2	7.2	7.2
IV_FreqSevB	6.8	5.9	3.3	3.9	3.9	3.9	11.0	8.5	-1.6	8.9	10.3	10.3	10.3	10.1	0.0	9.9	-1.6	-1.3	10.1	0.0	9.9	-1.6	-1.3	11.0	11.0
IV_FreqSevC	3.4	1.8	0.0	-0.2	-0.2	-0.2	7.4	3.9	-0.9	4.5	4.5	4.5	4.5	0.8	-1.7	0.0	-0.9	-0.9	0.8	-1.7	0.0	-0.9	-0.9	7.4	7.4
DepRatio1	6.8	6.3	3.2	4.1	4.1	4.1	11.3	9.0	-2.3	9.5	11.0	11.0	11.0	10.4	4.4	10.2	0.0	-0.5	10.4	4.4	10.2	0.0	-0.5	11.3	11.3
DepRatio36	6.8	6.2	3.2	4.0	4.0	4.0	11.2	8.9	-2.0	9.5	11.0	11.0	11.0	10.4	4.0	10.2	0.9	0.0	10.4	4.0	10.2	0.9	0.0	11.2	11.2

Table 5: Gini Standard Errors for Fourteen Scores

Base Premium	Basic Explanatory Variables				Extended Explanatory Variables									
	Single Peril		IND_		Single Peril		IND_		Pure		Pure		IV_	
	Freq	Prem	Freq	Prem	Freq	Prem	Freq	Prem	Freq	Prem	Freq	Prem	Freq	Prem
ConsPrem	1.43	1.43	1.43	1.42	1.41	1.45	1.43	1.34	1.43	1.34	1.39	1.34	1.38	1.45
SP_FreqSev_Basic	0.00	1.43	1.43	1.46	1.46	1.42	1.42	1.45	1.45	1.45	1.39	1.44	1.38	1.45
SP_PurePrem_Basic	1.44	0.00	1.41	1.41	1.45	1.45	1.42	1.46	1.46	1.46	1.42	1.45	1.40	1.47
IND_PurePrem_Basic	1.44	1.39	0.00	1.38	1.38	1.40	1.40	1.44	1.44	1.44	1.40	1.42	1.38	1.34
IV_PurePrem_Basic	1.45	1.43	1.39	0.00	1.45	1.45	1.48	1.41	1.41	1.41	1.40	1.40	1.39	1.46
SP_FreqSev	1.47	1.49	1.45	1.50	1.50	0.00	1.44	1.47	1.47	1.46	1.45	1.47	1.45	1.49
SP_PurePrem	1.42	1.43	1.42	1.49	1.49	1.41	0.00	1.45	1.45	1.45	1.46	1.46	1.46	1.49
IND_FreqSev	1.47	1.48	1.48	1.45	1.45	1.45	1.48	0.00	1.38	1.40	1.53	1.30	1.50	1.52
IND_PurePrem	1.43	1.45	1.33	1.44	1.44	1.47	1.48	1.49	1.49	1.49	1.50	1.50	1.49	0.00
IV_PurePrem	1.44	1.47	1.41	1.42	1.42	1.48	1.51	1.45	1.45	1.45	1.49	1.45	1.50	1.46
IV_FreqSevA	1.41	1.43	1.45	1.43	1.43	1.43	1.49	1.54	1.54	1.54	0.00	1.55	1.34	1.53
IV_FreqSevB	1.47	1.47	1.47	1.44	1.44	1.46	1.49	1.30	1.30	1.27	1.55	0.00	1.53	1.54
IV_FreqSevC	1.41	1.42	1.43	1.43	1.43	1.44	1.49	1.52	1.51	1.51	1.34	1.53	0.00	1.53
DepRatio1	1.47	1.48	1.48	1.45	1.45	1.45	1.48	1.38	0.00	1.32	1.53	1.30	1.50	1.52
DepRatio36	1.48	1.48	1.48	1.45	1.45	1.45	1.48	1.40	1.32	0.00	1.53	1.27	1.50	1.49

Table 6: Gini Indices, Approximations and Decompositions

Score	Gini	Gini_approx	LossSource	PremSource	GiniSource	Gini_approx2	GiniReverse	SumGinis
SP_PurePrem_Basic	4.89	4.83	-6.14	-10.97	-7.58	3.82	4.34	9.37
IND_PurePrem_Basic	4.01	3.98	-3.60	-4.45	-8.87	3.81	8.08	12.49
IV_PurePrem_Basic	4.33	4.42	9.25	-2.03	-5.53	4.91	7.88	12.75
SP_FreqSev	11.15	11.29	4.51	-9.72	-9.72	10.79	1.71	13.18
SP_PurePrem	9.97	10.04	0.59	-9.56	-7.16	9.15	6.97	16.92
IND_FreqSev	10.03	10.31	1.65	-7.80	-9.29	10.71	6.85	17.39
IND_PurePrem	10.96	11.21	4.28	-9.60	-9.69	10.10	6.13	17.60
IV_PurePrem	11.29	11.44	5.06	-9.69	-9.69	10.56	6.69	18.26
IV_FreqSevA	12.59	12.86	1.54	-9.60	-9.69	12.31	2.79	15.64
IV_FreqSevB	10.61	10.83	5.89	-9.60	-9.69	11.18	6.77	17.76
IV_FreqSevC	12.80	12.99	0.76	-9.60	-9.69	12.31	3.44	16.37
DepRatio1	10.09	10.36	0.65	-9.60	-9.69	10.74	6.83	17.41
DepRatio36	10.06	10.34	0.65	-9.60	-9.69	10.72	6.77	17.33

Note: Base Premium is *SP\_FreqSev\_Basic*.

Table 7: Gini Indices from a Simulation Study

Scenario	Base		Simple Gini	Ratio Gini Indices								True Correlation
	Premium	Score1		Score2	Score3	Score4	Score5	Score6	Score7	Score8	maximum	
Explanatory variables contribute equally	Score1	10.05	0.00	0.10	-0.09	0.06	-0.02	0.02	-0.12	-0.03	0.10	99.29
	Score2	7.13	7.08	0.00	4.97	-0.12	4.29	7.08	6.42	6.08	7.08	68.78
	Score3	7.08	7.13	5.12	0.00	4.38	0.00	6.61	7.13	6.25	7.13	68.77
	Score4	7.13	7.79	2.73	5.02	0.00	4.29	7.67	7.08	7.09	7.79	68.78
	Score5	7.08	7.75	5.11	2.64	4.32	0.00	7.13	7.59	7.13	7.75	68.77
	Score6	9.75	2.63	-1.23	1.08	-0.98	-0.09	0.00	1.87	-0.12	2.63	96.90
	Score7	9.70	2.73	1.38	-1.46	0.09	-1.04	2.19	0.00	0.02	2.73	96.69
	Score8	9.57	4.13	0.59	0.27	-1.24	-1.47	2.73	2.63	0.00	4.13	94.65
Second explanatory variable contributes little	Score1	7.19	0.00	0.40	-0.08	0.36	-0.01	0.30	-0.12	0.03	0.40	98.48
	Score2	7.13	1.44	0.00	0.34	-0.12	0.25	1.34	0.80	0.64	1.44	96.49
	Score3	1.44	7.13	7.04	0.00	6.67	0.31	7.07	6.96	6.88	7.13	21.33
	Score4	6.95	3.13	2.74	-0.48	0.00	-0.64	3.07	1.44	1.34	3.13	94.56
	Score5	1.30	7.16	7.10	0.50	6.78	0.00	7.13	7.01	6.96	7.16	20.38
	Score6	7.15	0.52	0.27	0.08	0.28	-0.08	0.00	0.22	-0.11	0.52	97.64
	Score7	6.93	2.74	2.67	-1.22	0.40	-1.14	2.67	0.00	0.29	2.74	95.79
	Score8	6.89	2.87	2.71	-1.02	0.27	-1.22	2.74	0.53	0.00	2.87	94.86
First explanatory variable contributes little	Score1	6.91	0.00	0.36	0.63	0.40	0.56	0.16	0.41	0.34	0.63	98.47
	Score2	0.52	6.94	0.00	7.00	0.36	6.92	6.94	6.94	6.90	7.00	5.35
	Score3	6.94	0.53	0.44	0.00	0.46	0.17	0.40	0.52	0.43	0.53	99.08
	Score4	0.52	6.95	0.25	7.00	0.00	6.91	6.91	6.94	6.94	7.00	5.85
	Score5	6.94	2.56	-0.75	2.51	-0.73	0.00	0.53	2.53	0.52	2.56	99.08
	Score6	6.82	2.53	-0.97	2.73	-0.87	0.62	0.00	2.57	0.41	2.73	97.53
	Score7	6.90	0.20	0.41	0.66	0.35	0.61	0.29	0.00	0.16	0.66	98.11
	Score8	6.81	2.52	-0.87	2.72	-0.98	0.67	0.21	2.53	0.00	2.72	97.05