

→ Thesaurus ! (✓) BUY: bought, investing, acquisition
useful.
2007-20 pages - 2021.06.23
- a bit outdated
- but still relevant
ResearchGate

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/228258309>

Market Reaction to Verbal Components of Earnings Press Releases: Event Study Using a Predictive Algorithm

Article in Journal of Emerging Technologies in Accounting · January 2007

DOI: 10.2308/jeta.2006.3.1.1

CITATIONS

101

READS

799

1 author:



Elaine Henry

Stevens Institute of Technology

66 PUBLICATIONS 1,373 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Quantitative Modeling of Unstructured Financial Information [View project](#)

Market Reaction to Verbal Components of Earnings Press Releases: Event Study Using a Predictive Algorithm

Elaine Henry
University of Miami

ABSTRACT: Similar to a classic-event study, this study examines market reaction to firms' earnings announcements. This study extends the examination to include a broad range of concurrent disclosure contained in earnings press releases: financial disclosure captured as accounting ratios; and verbal components of disclosure, both content and style, which are captured using elementary computer-based content analysis. Extending the analysis to such a broad range of concurrent disclosures requires a methodology designed to utilize a large number of predictor variables, and predictive data mining algorithms are specifically designed to do so. Therefore, this study employs a widely used data-mining algorithm—classification and regression trees (CART). Results of the study show that inclusion of predictor variables capturing verbal content and writing style of earnings-press releases results in more accurate predictions of market response.

Keywords: CART; content analysis; event study; earnings announcements.

INTRODUCTION

Press releases announcing periodic earnings, though voluntary, are an important means by which firms communicate with investors. Recent changes in regulation, as well as advances in technology, have increasingly made these press releases directly accessible to investors. Following the Sarbanes-Oxley Act of 2002, the Securities and Exchange Commission (SEC) now requires that all earnings press releases must be filed, thus ensuring that investors have direct access to the actual text of firms' earnings releases.¹ Even apart from this requirement, advances in technology such as broadcast fax and the internet have given investors direct access to these press releases. To gain a better understanding of aspects of this increasingly direct firm-to-investor communication process, this study explores market reaction to verbal components of firms' earnings press releases.

¹ As a result of the Sarbanes-Oxley Act of 2002 ("SOX"), the Securities and Exchange Commission (SEC) has added a new disclosure—Item 12 to Form 8-K, formally requiring firms to file all news releases announcing earnings (SEC, 2003).

I thank Glenn Shafer (dissertation committee chair), Elizabeth Gordon, Darius Palia, and Karen Kukich (National Science Foundation) for helpful suggestions on an earlier version of this paper, and the editor Miklos Vasarhelyi and anonymous reviewers for their constructive comments. I acknowledge, with gratitude, support from the Whitcomb Center for Research in Financial Services and financial support from Deloitte and Touche Foundation.

Corresponding author: Elaine Henry
Email: e.henry1@miami.edu

Earnings announcements and the market reaction to the information contained within them have been the subjects of a substantial body of empirical-capital-markets research. Research has shown that the relationship between periodic accounting results and event-window market returns, though significant, is small. Some accounting researchers have attributed the small magnitude of this relationship to irrelevance of accounting information (Lev 1989), while others have proposed that the small magnitude of this relationship is to be expected because accounting information reflects historical performance, and market price movements reflect changed expectations about future profitability (Kothari 2001).

Abrahamson and Amir (1996), in an analysis of annual report disclosure, suggest that another explanation for the low association between returns and earnings is the omission of textual components of disclosure from these analyses. This explanation can also be applied to earnings announcements, and is supported by accounting research showing that disclosures made concurrently with earnings announcements, including officer comments, help explain market reaction (Francis et al. 2002; Hoskin et al. 1986).

The purpose of this study is to examine whether concurrent disclosure, particularly information contained in the verbal components of earnings press releases and captured with keyword counts, improves prediction of market response to earnings announcements. Similar to a classic-event study, this study examines market reaction to summary measures of earnings performance, e.g., change in EPS. To include concurrent disclosure in the analysis, this study extends the number of predictor variables to capture additional aspects of the financial and verbal information disclosed in earnings press releases. The additional financial predictor variables include key accounting ratios, and the verbal predictor variables include keyword counts and other measures aimed at capturing both content and style of firms' commentary on their performance.

Including verbal components of the press releases in the analysis of market reaction requires a methodology to capture both content and style. Verbal content is captured using elementary computer-based content analysis to develop counts of occurrences in each press release of keywords (e.g., dividends, sales). To organize the keywords, a thesaurus relating specifically to periodic financial reporting is constructed based on the hand-coding schemes used by accounting researchers to categorize information in earnings press releases (Hoskin et al. 1986; Kasznik and Lev 1995; Miller 2002; Shon 2003), and organized loosely around a resource-enterprise-agent framework (McCarthy 1982). Style is captured using measures of tone, length, numerical intensity, and complexity as used in Henry (2005).

Including a broad range of measures aimed at capturing concurrent disclosures requires a methodology capable of utilizing a large number of predictor variables, and predictive-data-mining algorithms are specifically designed to do so. This study employs the widely used data-mining algorithm classification and regression tree (CART). Using the dependent variable of abnormal market returns around the announcement date (i.e., a binary variable 0/1) indicating negative or positive excess returns versus the market, this study examines the impact of the various types of predictor variables on predictive accuracy.

Results show that inclusion of a broad range of financial variables, such as levels and changes of ratios, does not enhance predictive accuracy of market returns incrementally above a model including simple earnings information. However, inclusion of verbal variables does result in greater predictive accuracy.

The remainder of the paper is organized as follows: The second section describes classification trees and the specific algorithm employed in this study. The third section provides background on predictor variables and their relationship to existing research. Following this section are sections that describe the sample and data collection method and present results. The final section summarizes and concludes.

CLASSIFICATION TREES

Data-mining approaches are becoming generally more important as a result of the rapid growth of information and databases. An exploration of this technique in capital-markets-accounting research is warranted simply because of its widespread use in other data-intensive disciplines. The tree approach, and Classification and Regression Trees (CART) in particular, is a widely used data-mining algorithm (Hand et al. 2001).

Classification tree algorithms (and data mining approaches in general) have the general property that high-dimensional data sets do not necessarily degrade the quality of the results. The ability of the techniques to take advantage of a large number of variables is particularly useful in the current study because of the large number of explanatory variables employed to capture the effect of concurrent disclosure on market reaction.

One implication of using a classification tree rather than ordinary least squares (OLS) regression, the customary method used in event studies, is that the research question is altered so that the goal is to make the best prediction from the data—not to investigate coefficients of specific variables. In this study, market reactions are classified as either over-performance or under-performance relative to the market on the day of the earnings announcement, and the research question is whether inclusion of variables capturing concurrent disclosure—both financial and verbal—improves prediction of market reaction.

The following sections include a description of the classification-tree approach and rationale for use of CART in this study.

Description of the Classification-Tree Approach

This study uses CART (specifically, a commercial product of the same name, CART, distributed by Salford Systems) to predict market reaction to both verbal and numerical components of press releases. The classification tree approach in CART, termed “binary recursive partitioning” (Breiman et al. 1984), involves repeatedly splitting the data set into two subgroups that are relatively homogenous with respect to the dependent variable. Each split is made by evaluating all possible splits and choosing a variable and a level of that variable that make the two resulting subgroups as homogeneous as possible. (See Hastie et al. [2001, 271] for discussion of the splitting criteria, known as the Gini index.) Because each split of a group can be graphically represented as creating two branches, this recursive splitting is also known as tree growing.

Based only on training data, a tree could be grown with zero classification error if, for example, there were one node for each observation. Such a tree may or may not perform well on other data, i.e., a problem of over-fitting. To estimate how well a tree will perform on new data, CART uses a process of cross-validation which repeatedly partitions the data into a training portion used to grow an auxiliary tree and a portion for testing. In ten-fold cross validation, 9/10 of the data is used to grow each auxiliary tree, and 1/10 of the data is used to test the predictive accuracy of that tree; this process is repeated ten times, and the test error rates are averaged to provide an estimate of the error rate for a final tree (which is based on all of the data). The process of cross-validation is also used in determining the optimal-sized tree, one that is complex enough to result in good predictions but not so complex that it would not perform well on new data. In summary, the best tree is that sub-tree with the lowest cross-validated test error rate. (See Breiman et al. [1984] for further explication of the CART methodology.)

Use of CART in this Study

A primary reason for using the CART methodology in this study, as noted above, is its general property that high-dimensional data sets do not degrade results—particularly

important because this study of concurrent disclosure uses a large number of predictor variables. In addition, the methodology has the ability to deal with missing data, outliers, independent non-recoded categorical variables, and nonlinear relationships. Classification trees, are not, however, causal models and do not generally give insight into causal relationships.

Although the research question in this paper is how inclusion of additional/verbal information improves prediction—not a comparison of the predictive ability of various tools—a discussion of differences between the tree approach and classical statistical approaches may be helpful because the tree approach is less common in accounting capital markets research.

Logistic regression is an approach often used with binary dependent variables such as that used in this study. Predicted classification of observations is based on the estimated probability exceeding one-half, and the goodness-of-fit of a logit model can be judged by totaling the proportion of each class correctly classified (Kennedy 1998). Although the application of CART in this study is similar to this use of logistic regression in terms of ultimate objective and in terms of evaluating a model's success, major differences exist.

The first major distinction between the tree approach and classical statistical approaches is that the tree approach is an algorithmic approach with roots in machine learning, while logistic regression (and ordinary least squares, OLS) are data model approaches with roots in classical statistics. In logistic regression and OLS, the goal is to find a model of the way things work, and in general, additional variables degrade results. In contrast, the classification-tree approach is aimed at prediction and has the general property that more variables are less likely to degrade the value of the prediction. As Breiman (2001, 199) writes: “There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown.”

As an algorithmic approach, the tree approach “learns” decision rules (for classification or regression) from the data, using a training set of data, and the results are judged by how well the decision rules perform on the test set of data. A logistic regression aims to estimate the parameters of the model. In contrast, the primary focus of CART is the predictive success of a model; and although the relative importance of various predictors is identified, there is no direct equivalent of a single coefficient β because a single predictor variable may be used in various levels of a tree. CART is a recursive algorithm, partitioning the data at each node on a single variable; in contrast, logistic regression involves all predictor variables simultaneously and uses maximum likelihood to estimate coefficients.

CART and logistic regression are two separate tools that can also be used in combination, for example, using CART to initially identify relevant break points in variables' values that can then be included in a logistic regression.

The large number of predictor variables used in this study also results in a data set with missing values and creates issues with the use of logistic regression, which would require a complete set of values for every observation. Missing data has been shown to be a problem in capital markets research. For example, Ou and Penman (1989a, 1989b), in their analysis of the ability of a large number of financial statement attributes to predict the direction of future earnings, cite missing data as a motivation for their multi-stage variable reduction process (a univariate logit for initial variable reduction, followed by a step-wise regression).

The problem of missing data is addressed by the tree classification approach. CART processes observations with missing data for some variables by developing surrogate splitting variables that are used when data on the primary splitting variable is missing (Hastie

et al. 2001). To develop surrogate-splitting variables, after identifying the best predictor variable and value for each split, CART identifies other variables and values that would have resulted in similar classification outcomes as the variable actually chosen for that split. CART ranks these surrogate variables on how similar the classification outcomes would have been in comparison with the variable actually chosen. Then, in instances where an observation is missing a value for a particular variable, CART selects the highest-ranked surrogate variable for which a value is available.

It may also be noted that the tree approach bears only a superficial resemblance to the variation of OLS using stepwise regression. Unlike parametric stepwise procedures, the tree approach is statistically sound. The tree approach uses only cross-validation in prediction, unlike stepwise procedures that use goodness-of-fit measures. Stepwise regression is done either by using minimum and maximum F -statistics for keeping or eliminating particular independent variables (Neter et al. 1996) or by sequentially regressing the dependent variable on each independent variable and keeping the model that produces the highest R^2 (Kennedy 1998). However, stepwise regression has statistical problems (particularly when the number of variables is large relative to the number of observations), including inflated R^2 , F , and t -statistics (Freedman 1983; Lovell 1983; Rencher and Pun 1980). As noted, owing to cross-validation, the tree approach does not face the same statistical shortcomings.

In summary, the strengths of the tree approach relevant to this study include most importantly, the property (common to data-mining algorithms generally) that high-dimensional data sets do not degrade results; and additionally, the ability to deal with missing data, outliers, independent categorical variables, and nonlinear relationships. Overall, despite certain weaknesses, the strengths of the CART methodology make it well suited to the objective of examining the market impact of a broad range of predictor variables capturing concurrent financial and verbal disclosure.

PREDICTOR VARIABLES CAPTURING CONCURRENT DISCLOSURE

Using the dependent variable of abnormal market returns around the announcement date (i.e., a binary variable indicating negative or positive excess returns versus the value-weighted S&P 500 returns on the day of earnings announcement), this study examines five types of predictor variables: (1) longer-term firm characteristics, such as industry and size; (2) current period earnings information; (3) levels and changes of other financial information; (4) keyword counts capturing the nature of operating information disclosed in the verbal components of press releases; and (5) style, including length of press release, complexity, numerical intensity, and tone. The following paragraphs describe each type of predictor variable, along with references relating each to market returns.

The first type of predictor variable captures longer-term firm characteristics. These predictor variables, collectively referred to as *FIRM-Info* type include:

- Size (Market Value *MKVAL*, Total Assets *AT*). Abnormal returns during earnings announcement periods are related to firm size (Ball and Kothari 1991).
- Previous year unexpected earnings *Lag_UE*. Firms reporting unexpectedly high earnings experience superior market performance in the following periods (Ball and Brown 1968; Bernard and Thomas 1989; Foster et al. 1984).
- Market performance over previous 48 trading days (i.e., roughly since the last earnings announcement) *CMXSPREV48*. Return-based trading strategies, whether momentum or contrarian (DeBondt and Thaler 1985; Jegadeesh and Titman 1993) are related to autocorrelation of stock prices. Since the autocorrelation of stock prices is a distinct phenomenon from the market impact of prior period's earnings surprises

(Chan et al. 1996), this predictor variable captures market performance roughly since the last earnings announcement.

- String of previous years (within the past five) with positive net income *POSNI*, string of previous years (within the past five) with positive *EBITDA* (earnings before interest, taxes, and depreciation) *POSEBITDA*. Burgstahler and Dichev (1997) show that firms' earnings management is increasing in the length of the preceding string of positive earnings and that the magnitude of earnings management to avoid economic losses is economically significant relative to firms' market value.
- Ownership (percentage institutional ownership *IOTSHR*, and number of common shareholders *CSHR*). Prior research shows evidence that abnormal returns around earnings announcements are negatively related to the percentage of institutional ownership (Bartov et al. 2000; Hotchkiss and Strickland 2003).
- Industry (*SIC*). Although Moskowitz and Grinblatt (1999) cite evidence that industry has not been found to be important in explaining asset prices generally, they provide evidence that industry is important in explaining results of momentum-trading strategies.

The second type of predictor variable, collectively referred to as the *EARN-Info* type, captures current earnings information. Aspects of the earnings information include changes in performance from the previous year, performance relative to analysts' expectations, and whether the firm reported a profit or a loss. These variables are as follows:

- Earnings disclosed in the current press release compared to previous year earnings (*SRW_UE*). Following the approach in Francis et al. (2002), unexpected earnings are measured as actual minus the previous year, scaled by share price.
- Earnings compared to analysts' forecasts (indicator equal to 1 if company's reported EPS exceeded the average of analysts' forecast *DIF_gt0*). Market reaction to positive (versus negative) unexpected earnings relative to analysts' forecasts differs (Lopez and Rees 2002).
- Earnings positive or negative (indicator equal to 1 if company reported a profit in the current year, 0 if a loss *NI_gt0*). Hayn (1995) shows that the market response to firms with profits is higher than for those with losses.

The third type of predictor variable captures other financial information disclosed in earnings announcements, including levels and changes in accounts and in commonly used ratios involving these items, collectively referred to as the *OTHER-FIN* variable type. This list extends that used by Ou and Penman (1989a, 1989b) in their step-wise regression analysis predicting future accounting results from reported accounting results. These variables can be derived from information contained in the full-earnings results, although not all of the items are always disclosed in a company's earnings' press release. Appendix A lists these variables.

The fourth type of predictor variable, collectively referred to as the *KEY-WORDS* variable type, uses keyword counts to capture the nature of the information discussed in the verbal components of the earnings press releases. The keyword counts were obtained using commercially available software (Diction 5.0 available at: <http://www.dictionsoftware.com/>). The specific software was selected for the following reasons: (1) it is Windows-based and easily learned; (2) it has the capability to process multiple texts and produce summary statistics; and (3) it allows custom thesauruses in computing term frequencies.

Term-frequency count metrics have been used in research on accounting narratives usually in annual reports. For example, Smith and Taffler (2000) show that particular combinations of words and phrases in annual reports are associated with bankruptcy. Hussainey et al. (2003) use keywords associated with forecasts and predictions to quantify the number of sentences in annual reports dealing with forward-looking statements, which increase the relationship between current stock returns and future earnings changes.

Term-frequency counts are also widely used in other domains. Using term-frequency counts to capture the content of a document is commonly used for textual analysis in the social sciences, where such techniques fall into the general category of analysis referred to as content analysis (Neuendorf 2002). Term-frequency counts, also described as treating a document as a “bag of words,” are also commonly used both for text-based information retrieval and for text categorization (Hand et al. 2001; Manning and Schutze 1999). Information retrieval and text categorization are two applications that fall into the overall category of automated language analysis, referred to in the computer sciences as natural language processing (NLP). Hand et al. (2001, 457) note:

NLP techniques (which try to explicitly model and extract the semantic content of a document) are not the mainstay of most practical [text-based information retrieval] systems in use today; i.e., practical retrieval systems do not typically contain an explicit model of the meaning of a document ... and instead rely on simple term matching and counting techniques, where the content of a document is implicitly and approximately captured (at least in theory) by a vector of term-occurrence counts.

As a basis for the keyword counts capturing the information discussed in the verbal components of earnings press releases, this study constructs a thesaurus of keywords relating specifically to periodic financial reporting. The thesaurus is based on the hand coding schemes used by accounting researchers to categorize information in press releases (Hoskin et al. 1986; Kasznik and Lev 1995; Miller 2002; Shon 2003), and organized loosely around the resource-enterprise-agent framework (McCarthy 1982). Appendix B presents this thesaurus.

The use of a topic-specific thesaurus in the study reported here contrasts with other approaches to computer-based verbal analysis in accounting research. Some approaches rely on subject-matter knowledge and case-by-case coding judgment. For example, Clatworthy and Jones (2003) use computers to obtain frequency counts, but semantic coding is done manually. Their study shows equivalent emphasis on positive aspects of performance in annual-report-accounting narratives for both good and poor-performing companies. Other approaches rely not on subject-matter knowledge to capture the information in text, but rather on patterns in the text. For example, Frazier et al. (1984) use factor analysis to develop themes in management’s annual report discussion, which are then used as independent variables in a discriminant model to classify firms by positive or negative returns in the following year. As another example, Kloptchenko et al. (2004) encode every word and every sentence from three companies’ quarterly reports, use Euclidian distances to capture similarities across documents, create clusters of similar reports, and compare membership in report clusters with the financial performance of the companies over the same time period.

Hand et al. (2001) contrasts the thesaurus approach with latent semantic indexing (an application of principal component analysis), both applicable to the issue of synonymy in term-count-frequency text representation. The thesaurus approach “uses a knowledge base (a thesaurus or ontology) that is created a priori with the purpose of linking semantically related terms together” (Hand et al. 2001, 465). Use of a topic-specific thesaurus in the study reported here is applicable for two main reasons: First, earnings press releases may

be considered a sub-language, i.e., with a limited domain of discourse; and second, previous accounting researchers have constructed detailed hand-coding schemes designed to capture the content of earnings announcements.

Using this topic-specific thesaurus, this study obtains document-by-document frequency counts of each set of words. These frequency counts are scaled by the total number of words in the text portion of the press release. In summary, each press release is represented as a vector of scaled frequency counts of the keywords, grouped according to the thesaurus.

The fifth type of predictor variable relates to writing style: length, tone, textual complexity, and numerical intensity. These variables are collectively referred to as the *STYLE* type of variables:

- Length is measured as word count of the total press release *WDS_AL*, the commentary portion *WORDS*, and the financial statement portion *WDS_FIN*. Length, or quantity, is often used to measure disclosure quality on the grounds that regulation and litigation threats, as well as reputation effects, will assure that all disclosure is of high quality (Botosan 1997).
- Tone is measured as the frequency count of positive and negative words, scaled by total words. The words counted as positive and negative within earnings press releases are included in Appendix B. Positive tone (*POSITIVITY*) and negative tone (*NEGATIVITY*) are defined as the frequency counts of words on these lists, scaled by the total word count of the press release. Henry (2005) shows a relationship between tone as measured in this way and investor reaction to earnings announcements. Related research on specific words in accounting narrative—chairman’s letters rather than earnings announcements—has shown evidence of a relation between negative words in the letter and market returns (Abrahamson and Amir 1996) and firm failure (Smith and Taffler 2000).
- Textual complexity *COMPLEXITY*, or basic readability, is a widely used measure of style. Securities and Exchange Commission (SEC) regulations require “plain English” in disclosure (SEC 1998), and textual complexity is one of the most basic measures in content analysis (Neuendorf 2002). This study uses the measure of textual complexity provided in Diction 5.0, number of characters per word.
- Numerical intensity *NUMERICAL* refers to the amount of quantitative information in financial disclosure. Quantitative information is considered to be of higher quality and more credibility-enhancing than non-quantitative because of its precision Botosan (1997). This study measures numerical intensity as the percent of the textual portions, i.e., not including the financial statements, of a press release that is numerical: numerical terms² divided by word count.

SAMPLE AND DATA COLLECTION METHOD

The sample includes data on 441 companies for the year 2002, from the telecommunications industry, computer services industry, and related equipment manufacturers. During this period, close, notable linkages existed between telecommunications and computer services companies, implying a similar disclosure environment.³ The main rationale for using

² Numerical terms include integers, numbers in lexical format such as “one,” and terms referring to numerical operations (Hart 2000).

³ Fama and French (1997) define telecommunications as SIC codes 4800-4899; these SIC codes cover telecommunications, radio, television, and cable. Related equipment manufacturers include SIC codes 3661-3669. The computer-services industries include SIC codes 7370-7374 (programming, software, systems design, processing). Related equipment manufacturers include SIC codes 3570-3579. The source for the SIC code list is <http://www.sec.gov/>.

this industry and time period is the greater uncertainty created by the stock market boom of the late twentieth century. Research suggests that investors rely more heavily on nonfinancial information in periods of uncertainty and rapid technology change (Amir and Lev 1996). Further, the industry offers varied performance.

Companies included in the sample are those with a fiscal year-end of December, share price greater than \$1, and whose information matched between CRSP and Compustat databases. Accounting variables are obtained from the Compustat database, return variables from the CRSP database, and analyst forecast information from the IBES database.

Earnings press releases are obtained from Nexis-Lexis and Factiva. The press releases are those issued directly by the company and are identified by source (typically on *PRNews-wire* or *Business Wire*), by length (longer than non-company authored reports), and ultimately by inclusion of a company contact (usually an investor relations contact). Each press release is saved as a separate text file, with an identifier included to allow for subsequent remerging of term-frequency count data with other data. Rather than stemming, an approach which transforms words into a basic stem of the word, the thesaurus used in this study includes various forms of the word. For example, both “revenue” and “revenues” are included in the thesaurus.

From the earnings press releases, the content-analysis software is employed to obtain document-by-document term-frequency counts organized according to the topic-specific thesaurus described above. The term-frequency counts are scaled by the total number of words in the text portion of the press release. Overall, the approach taken in this study captures content as a collection of synonym sets. Each press release is represented as a vector of scaled-frequency counts of keywords, grouped according to the thesaurus.

RESULTS

This section reports the results of using various types of explanatory variables to predict abnormal market returns around the announcement date. Prediction success rate is first defined and then used to compare the predictive accuracy of classification trees using different sets of information.

The prediction success rate is simply 1 minus the misclassification rate. The rule used in CART for selecting the optimal tree (minimization of misclassification errors) is expressed in terms of misclassification costs. Because this study assumes that it is no more costly to misclassify a class 1 case (positive excess market returns) than a class 0 case (negative excess market returns), results expressed as misclassification costs are no different than simply the proportion of objects misclassified, i.e., the misclassification rate. Another parameter, which is assumed in tree classification, is the prior probability of an observation being in a class. The current study assumes that each case had a 50/50 chance of being in class 1 or class 0, i.e., underperforming or outperforming the market on earnings-announcement day.

An example of prediction success is provided in Table 1. The predictions were made based on information set *FIRM-Info*, including only the longer-term firm characteristics. As shown, 122 (55.71 percent) of the total 219 class 0 observations were correctly classified as class 0. Of the total 222 class 1 observations, 127 (57.21 percent) were correctly classified. The overall success rate is therefore 56.46 percent.⁴

Table 2 summarizes the cross-validated test sample results for each information set included, which provide an out-of-sample estimate of performance. Results are presented

⁴ 56.46% = 55.71%*50% + 57.21%*50%.

TABLE 1
Classification Tree Example Results
Cross-Validated Test Sample Prediction Success and Misclassification

Actual Class	Total Cases	Percent Correct	Classed as 0^a n = 217	Classed as 1^b n = 224
0	219	55.71	122	97
1	222	57.21	95	127

^a Class 0 refers to under-performance of the market on earnings announcement day.

^b Class 1 refers to out-performance of the market.

These predictions are based only on information set 1, *FIRM-info*, containing longer-term firm characteristics. Of the total 219 class 0 observations, 122 observations (55.71 percent) were correctly classified as class 0. Of the total 222 class 1 observations, 127 (57.21 percent) were correctly classified.

TABLE 2
Cross-Validated Test Sample Prediction Success Rate for Each Information Set Included in the Classification Tree

Successful classification as either under-performance or out-performance of market on earnings announcement date

Information Set Included	Prediction Success Rate (Test by 10-fold Cross-Validation)
1. <i>FIRM-Info</i> (9 variables)	0.5646
2. <i>FIRM-Info</i> & <i>EARN-Info</i> (12 variables) ^a	0.5533
3. <i>FIRM-Info</i> & <i>EARN-Info</i> & <i>OTHER-FIN</i> (134 variables)	0.5412
4. <i>FIRM-Info</i> & <i>EARN-Info</i> & <i>OTHER-FIN</i> & <i>KEY-WORDS</i> (162 variables)	0.5814
5. <i>FIRM-Info</i> & <i>EARN-Info</i> & <i>OTHER-FIN</i> & <i>KEY-WORDS</i> & <i>STYLE</i> (169 variables)	0.5952
6. <i>FIRM-Info</i> & <i>KEY-WORDS</i> & <i>STYLE</i> (44 variables)	0.5908

^a When only firms which reported a profit are included, prediction success rate is .6205. Results from restricting the sample in this way for the other information sets (not tabulated) do not exhibit the same improvement in prediction.

The types of predictor variables include:

FIRM-Info = longer term firm characteristics;

EARN-Info = earnings information for the current period;

OTHER-FIN = other financial information;

KEY-WORDS = keyword counts capturing the nature of the operating information disclosed in the verbal components of press releases; and

STYLE = style variables, including length of press release, complexity, numerical intensity and tone.

only for the test sample. Training sample results (not shown) are, of course, better than test sample results.

The types of predictor variables (described above) combined to form the information sets include:

FIRM-Info = longer term firm characteristics;

EARN-Info = earnings information for the current period;

OTHER-FIN = other financial information;

KEY-WORDS = keyword counts capturing the nature of the operating information disclosed in the verbal components of press releases; and

STYLE = style variables, including length of press release, complexity, numerical intensity and tone.

These various types of predictor variables are combined to form different information sets, and then input as predictor variables in the CART data-mining algorithm. The primary research focus is on whether inclusion of verbal variables (both keyword counts and style measures) increases predictive accuracy.

As shown, Information Set 1 containing only variables measuring long-term firm characteristics *FIRM-Info* resulted in a cross-validated prediction success rate of 56.46 percent. Predictive accuracy using Information Set 2, which additionally includes contemporaneous earnings information *EARN-Info*, was lower at 55.33 percent (although when the sample is restricted to include only profitable firms inclusion of the earnings variables—substantially improves the prediction success rate).

Information Set 3 includes *FIRM-Info*, *EARN-Info*, *OTHER-FIN* (longer-term firm characteristics, earnings information and other financial variables) for a total of 134 predictor variables. As shown, the predictive accuracy of 54.12 percent with all financial variables is worse than results with the smaller information sets.

Information Set 4 adds keyword count variables *KEY-WORDS* for a total of 162 predictor variables, and results in an improvement of predictive accuracy, to a 58.14 percent prediction success rate. The prediction success rate of Information Set 5, which additionally includes the verbal-style variables *STYLE*, is 59.52 percent—better than the results with any of the other information sets.

Finally, Information Set 6 includes *only* longer-term firm characteristics *FIRM-Info* and the verbal variables (*KEY-WORDS* and *STYLE*). The relative prediction success rate of 59.08 percent is only slightly worse than the model with all earnings and financial variables included.

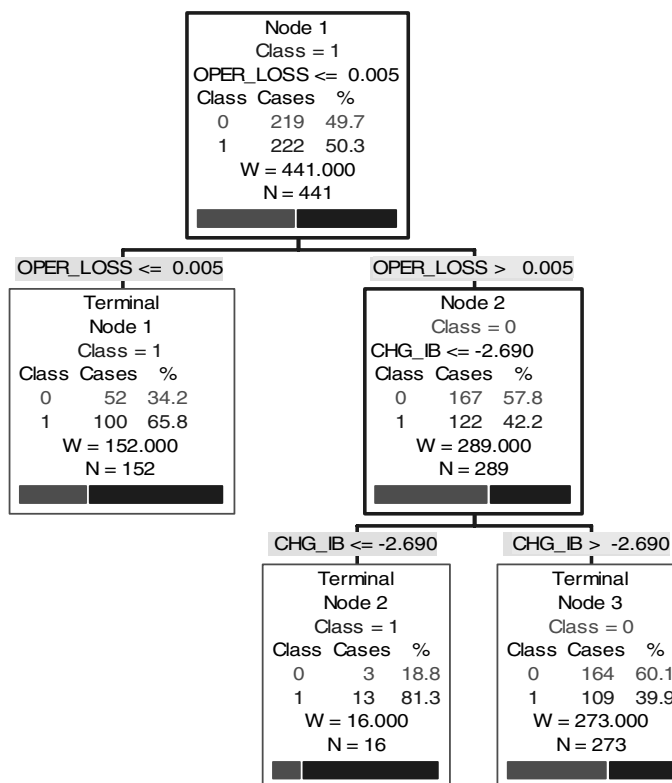
Comparing prediction success using Information Set 5, which includes all predictor variables, to the prediction success using Information Set 3, which includes all financial predictor variables but no verbal predictor variables, shows an improvement in prediction success of 5.40 percent.

Is this improvement meaningful? One way to evaluate this difference in estimated prediction success is to compare the improvement to the estimated standard errors of these success rates. Since this is a two-class problem, the count X of out-performers correctly identified would be binomially distributed, with mean np and variance $np(1-p)$. The proportion of correctly identified out-performers (count of out-performers divided by the number of observations, X/n) would have a mean equal to the mean of X divided by n , or $np/n = p$, and a variance equal to the variance of X divided by n^2 , or $(np(1-p))/n^2 = (p(1-p))/n$. If we consider the conservative case that the probability of misclassifying any single observation is 50 percent, given the sample size of $n = 441$ used in this study, the variance of the proportion of observations correctly classified is 0.000567, and the standard deviation is .0238, or 2.38 percent. Using this standard deviation as a conservative benchmark implies that the 5.40 percent improvement in prediction (which is 2.3 times the standard deviation) resulting from use of the expanded information set is a significant improvement.

Example Classification Tree

Figure 1 shows an example CART tree for Information Set 4, grown with one additional level beyond the minimal cost tree. In this tree, the first split is made based on the variable *OPER_LOSS* (count of words relating to losses and expenses, scaled by the total words in

FIGURE 1
Example Classification Tree
CART Tree for Information Set 4 Grown with One Additional Level Beyond the Minimal Cost Tree



Information Set 4 includes: *FIRM-Info*, *EARN-Info*, *OTHER-FIN*, *KEY-WORDS* (longer-term firm characteristics, earnings information, other financial variables, and keyword count variables). Tree shows training sample results.

- In this tree, the first split is made based on the variable *OPER_LOSS* (count of words relating to losses and expenses, scaled by the total words in the text portion of the press release). Observations with *OPER_LOSS* less than or equal to 0.005 were split into the node on the left (Terminal Node #1) and classified as Class 1 (i.e., correctly classified). Interestingly, it is the verbal variable *OPER_LOSS* rather than any of the financial measures upon which the first split is made and which results in the minimal cost tree.
- The second split is made according to *CHG_IB* (a financial variable measuring the percent change in income before extraordinary items). Those observations that had *OPER_LOSS* greater than 0.005 and *CHG_IB* less than or equal to -2.69 were split into the node on the left (Terminal Node #2) and classified as Class 1, i.e., market out-performance. Of the observations in this node, 81.3 percent were actually Class 1 (i.e., correctly classified). Among those companies with relatively more mentions of losses and expenses, the change in income before extraordinary items was an important variable for predicting market out-performance.

the text portion of the press release). Observations with *OPER_LOSS* less than or equal to 0.005 were split into the node on the left (Terminal Node #1) and classified as Class 1, i.e., market performance predicted to be above the S&P 500 on earnings announcement

date. Of the observations in this node, 65.8 percent were actually Class 1 (i.e., correctly identified).

Interestingly, it is the verbal variable *OPER_LOSS*, rather than any of the financial measures, upon which the first split is made and which results in the minimal-cost tree. Fewer mentions of losses and expenses in a company's earnings press release were an important predictor of market out-performance.

The second split is made according to *CHG_IB* (a financial variable measuring the percent change in income before extraordinary items). Those observations that had *OPER_LOSS* greater than 0.005 and *CHG_IB* less than or equal to -2.69 were split into the node on the left (Terminal Node #2) and classified as Class 1 (i.e., market performance predicted to be above the S&P 500 on earnings announcement date). Of the observations in this node, 81.3 percent were actually Class 1 (i.e., correctly classified). The interpretation of this node is that among those companies with relatively more mentions of losses and expenses, the change in income before extraordinary items was an important variable for predicting market out-performance.

CONCLUSION

In the study reported here, inclusion of verbal predictor variables improves performance of tree-based algorithms predicting over or under-performance of the market on earnings announcement day. When the predictor variables include measures to capture the content and style of the verbal component of earnings press releases, in addition to all financial information, predictive accuracy improves significantly—from around 54 percent to nearly 60 percent correctly classified.

Overall, the results of this classification-tree analysis provide evidence of market reaction to verbal components of earnings press releases: inclusion of verbal variables results in greater predictive accuracy beyond inclusion of financial information. A possible explanation for the greater predictive accuracy is that the information in the verbal components of earnings press releases is less stale than the numerical information. To a large extent, the numerical information is already public, if only indirectly, through the company's and analysts' forecasts and through information about the economy and industry as a whole. If the financial information in press releases is largely stale, the new information contained in earnings announcements may be largely in the verbal component—namely, what aspects of past and future performance a company will choose to emphasize in its commentary and what writing style, including tone, will be employed to communicate those results.

Results show that inclusion of financial variables such as levels and changes of ratios does not enhance predictive accuracy of market returns incrementally above a model including general firm information and simple earnings information. A possible explanation is that the earnings information is a parsimonious set of data that effectively summarizes the information in the broader range of financial ratios. Another, not necessarily exclusive, explanation is that market expectations are based on these items of earnings information (improved profitability, profit versus loss, and exceeding analysts' expectations) and thus they are more useful predictors of market reaction than the broader range of financial ratios. Another explanation is that predictive accuracy of the financial variables is reduced by the amount of redundant information. Further research will be helpful in exploring alternative explanations.

The results reported here are based on observations for firms in one industry, for one year, and hence the results may not generalize to other industries and/or time periods. This issue could be addressed by applying similar methodology to different samples. One might

conjecture that a sample taken from an industry or time period with less uncertainty might show less effect of information such as that contained in the verbal variables. On the other hand, given the upheaval faced by this industry during this time period and the resultant challenges to the credibility of disclosure, it is conceivable that observations from another industry and time period might show even greater effect of this type of information. This question is left for future study.

In general, constraints on analyzing the market influence of disclosure made concurrently with earnings announcements include: methodology, because large numbers of predictor variables reduce the power of statistical tests in commonly employed methods such as ordinary least squares; missing data, a relatively frequent occurrence in capital-markets research employing combinations of various databases; and labor intensity and subjectivity of hand-coding the verbal portion of earnings press releases containing firms' commentary on their performance. This study employs two tools to address these constraints: first, a predictive algorithm (CART), which is designed to utilize a large number of predictor variables and to be less affected by missing data; and second, elementary computer-based content analysis combined with a topic-specific thesaurus to code firms' verbal disclosure. Although computer coding of textual information may not be practical for many situations in accounting research, the volume of verbal financial disclosure and the labor intensity of analyzing this disclosure suggest clear benefits to use of computer tools for verbal analysis: decreased labor intensity of verbal data collection and consistency of coding. Further examination of the relative merits of alternative applications to achieve specific research objectives, as well as the sensitivity of prediction to alternative methodologies for capturing verbal information, will be useful.

Overall, the importance of data mining and text-processing applications can be expected to become increasingly important in accounting research with the continued increase of both numerical and natural language disclosure.

APPENDIX A OTHER FINANCIAL PREDICTOR VARIABLES

Variable Name and Description

<i>ABSXI_DV_SLS</i>	= Sum of absolute value of extraordinary items for 2000–2002 divided by sum of sales for 2000–2002.
<i>ADSLS02</i>	= Advertising Expense/Sales*
<i>ASRW_UE</i>	= Annual seasonal random walk unexpected earnings
Unexpected Earnings	= actual EPS diluted, excluding extraordinary items minus EPS for the previous year, scaled by beginning period price.
<i>CAPXY02</i>	= Capital Expenditures*
<i>CFLY02</i>	= Cash Flow*
<i>CFODT02</i>	= Cash Flow/Total Debt*
<i>CHEY02</i>	= Cash*
<i>PERCENT CHANGE</i>	= Assets
<i>COLLECTY02</i>	= Average Collection Period*
<i>CRY02</i>	= Current Ratio*
<i>CSHOY02</i>	= Common Shares Outstanding*
<i>DCEY02</i>	= LT Debt/Common Equity*
<i>DIF</i>	= Difference from expectation, calculated as actual EPS minus mean analysts' forecast, divided by the absolute value of the actual.
<i>DIVFL02</i>	= Dividends as percent Cash Flows*
<i>DLTISY02</i>	= Issuance of LT Debt*

<i>DLTRY02</i>	= Reduction in LT Debt*
<i>DPPPE02</i>	= Depreciation/Plant Assets*
<i>DTY02</i>	= Debt – Total*
<i>DTEQY02</i>	= Total Debt/Total Equity*
<i>DVPORY02</i>	= Cash Dividends-Common*
<i>DVY02</i>	= Dividend Payout*
<i>EMPY02</i>	= Employees*
<i>FREECFLY02</i>	= Free Cash Flow*
<i>GPMY02</i>	= Gross Profit Margin*
<i>IBY02</i>	= Income before Extraordinary Items*
<i>ICY02</i>	= Interest Coverage After Tax*
<i>INTANY02</i>	= Intangibles*
<i>INVTY02</i>	= Inventories-Total*
<i>INVXY02</i>	= Inventory Turnover*
<i>ISSDT02</i>	= Issuance Lt Debt/Total Debt*
<i>IVDAT02</i>	= Inventory/Total Assets*
<i>NIY02</i>	= Net Income (Loss)*
<i>NIFL02</i>	= Net Income/Cash Flow*
<i>NPMY02</i>	= Net Profit Margin*
<i>OIADPY02</i>	= Op Income After Depreciation*
<i>OIBDPY02</i>	= Op Income Bef Depreciation*
<i>PPENTY02</i>	= PP&E-Total Net*
<i>PPMY02</i>	= Pretax Profit Margin*
<i>PR02</i>	= Price at calendar year end 2002.
<i>PREQ02</i>	= Purchase of Treasury Stock as percent of Stock*
<i>PRSTKY02</i>	= Purchase Com & Pref Stock*
<i>QRY02</i>	= Quick Ratio*
<i>RDSLS02</i>	= R&D/Sales*
<i>RECTY02</i>	= Receivables-Total*
<i>RECXY02</i>	= Receivables Turnover*
<i>REDDT02</i>	= Reduction LT Debt/Total Debt*
<i>REVERSAL</i>	= 1 if change in sales for 2002 vs. 2001 was negative, & change in sales 2001 vs 2000 was positive.
<i>ROAY02</i>	= Return on Assets*
<i>ROEY02</i>	= Return on Equity*
<i>SALEY02</i>	= Sales-Net*
<i>SELLINVY02</i>	= Days to Sell Inventory*
<i>SEQY02</i>	= Stockholders' Equity*
<i>SLOWING</i>	= 1 if the growth in sales for 2002 vs. 2001 is less than the growth in sales for 2001 vs. 2000.
<i>SLSAR02</i>	= Sales/Accounts Receivable*
<i>SLSAT02</i>	= Sales/Total Assets*
<i>SLSCH02</i>	= Sales/Total Cash*
<i>SLSFA02</i>	= Sales/Fixed Assets*
<i>SLSIN02</i>	= Sales/Inventory*
<i>SLSWC02</i>	= Sales/Working Capital*
<i>SUE</i>	= Standardized Unexpected Earnings = actual EPS minus mean analysts' forecast, divided by the standard deviation of the estimates.
<i>WCAPY02</i>	= Working Capital*

WCAT02 = Working Capital/Total Assets*
XADY02 = Advertising Expense*
XDY02 = Depreciation-Amortization*
XRDY02 = R&D Expense*
ZSCOREY02 = Z Score-Measure of Bankruptcy*

* Also included percent change in this metric

APPENDIX B **Thesaurus Used as a Basis for the Keyword Count Variables**

Agents

ACCOUNTANTS: account accounts accounted accountant accountants accounting auditor auditors audit audits auditing audited

BANKERS: banker bankers creditor creditors

BOARD: board director directors chairman

COMPANY: company firm

INVESTORS: shareholder shareholders shareowner shareowners investor investors stockholder stockholders

LAWYER: law laws litigate litigates litigating litigated litigious litigation suit lawyer lawyers legal counsel attorney attorneys

MANAGER: manager managers management executive executives officer officers president employee employees

PRONOUNS: we our ours ourselves us I me mine my myself

Resources

ASSETS: asset assets cash receivable receivables inventory inventories property properties plant plants factory factories equipment goodwill research investment investments subsidiary subsidiaries venture ventures

DIVIDENDS: dividend dividends

EQUITY: equity shares common preferred stock buyback repurchase split

LIABILITIES: credit debt bonds borrow borrows borrowed borrowing loan loans liability liabilities interest coupon arrears bankrupt bankruptcy covenant covenants liquid liquidity illiquid illiquidity solvent solvency insolvent insolvency

Enterprises

AFFIRM: affirm affirms affirmed affirming confirm confirms confirmed confirming reaffirm reaffirms reaffirmed reaffirming reiterate reiterates reiterated reiterating

BUY: buy buys buying bought acquire acquires acquiring acquired acquisition acquisitions purchase purchases purchasing purchased invest invests investing invested investment investments

CUT: reduce reduces reducing reduced reduction cut cuts cutting eliminate eliminates eliminating eliminated

EXTERNAL: industry environment economy economic fundamentals market markets

FINANCING: list lists listed financing finance finances financed issue issues issuing issued offer offers offering offered shelf registration

FIRE: fire fires firing fired terminate terminates terminating terminated termination resign resigns resigning resigned resignation

FUTURE: future long-term longer-term forward-looking expect expects expecting expected expectations anticipate anticipates anticipated anticipating forecast forecasts forecasting forecasted continue continues continuing continued guidance outlook

HIRE: hire hires hiring hired appoint appoints appointing appointed appointment appointments

OPER_LOSS: loss lose loses losses lost deficit expense expenses expensive overhead overheads cost costs costly

OPER_PROF: sale sales revenue revenues income profit profits profitable profitability earn earns earnings earned returns gain gains ebitda margin margins shipments orders

RESTRUCTURE: restructure restructures restructured restructuring recapitalize recapitalization recapitalized recapitalizing recapitalizes

REVISE_RESTATE: revise revises revised revision revisions restate restates restated restatement restatements

SELL: sell sells selling sold dispose disposes disposing disposed disposal disposals divest divests divested divesting divestiture spins spin spinoff

STRATEGY: strategy strategic strategies strategically

TENDER: tender tenders tendered tendering

WRITEOFF: write writes wrote writing write-off writeoff write-offs impair impaired impairment

Tone

NEGATIVITY: disappoint disappoints disappointing disappointed disappointment risk risks risky threat threats threaten threatened threatening penalty penalties negative negatives negatively fail fails failed failing failure weak weakness weaknesses weaken weakens weakening weakened difficult difficulty hurdle hurdles obstacle obstacles slump slumps slumping slumped uncertain uncertainty uncertainties unsettled unfavorable downturn depressed down decrease decreases decreasing decreased decline declines declining declined fall falls falling fell fallen drop drops dropping dropped deteriorate deteriorates deteriorating deteriorated worsen worsens worsening worse worst low lower lowest less least smaller smallest shrink shrinks shrinking shrunk below under challenge challenges challenging challenged poor poorly

POSITIVITY: pleased delighted reward rewards rewarding rewarded opportunity opportunities enjoy enjoys enjoying enjoyed encouraged encouraging positive positives success successes successful successfully succeed succeeds succeeding succeeded accomplish accomplishes accomplishing accomplished accomplishment accomplishments strong strength strengths certain certainty definite solid excellent stellar good leading achieve achieves achieved achieving achievement achievements progress progressing deliver delivers delivered delivering leader leading up increase increases increasing increased rise rises rising rose risen double doubled doubles improve improves improving improved improvement improvements enhance enhances enhanced enhancing enhancement enhancements strengthen strengthens strengthening strengthened stronger strongest strongly better best more most above record high higher highest greater greatest larger largest grow grows growing grew grown growth expand expands expanding expanded expansion exceed exceeds exceeded exceeding beat beats beating

REFERENCES

Abrahamson, E., and E. Amir. 1996. The information content of the president's letter to shareholders. *Journal of Business, Finance & Accounting* 23 (8): 1157–1182.

- Amir, E., and B. Lev. 1996. Value-relevance of non-financial information: The wireless communications industry. *Journal of Accounting Economics* 22: 3–30.
- Ball, R., and P. Brown. 1968. An empirical evaluation of accounting income numbers. *Journal of Accounting Research* 6 (2): 159–178.
- , and S. Kothari. 1991. Security returns around earnings announcements. *The Accounting Review* 66 (4): 718–738.
- Bartov, E., S. Radhakrishna, and I. Krinsky. 2000. Investor sophistication and patterns in stock returns after earnings announcements. *The Accounting Review* 75 (1): 43–63.
- Bernard, V. L., and J. K. Thomas. 1989. Post-earnings-announcement drift: Delayed price response or risk premium? *Journal of Accounting Research* 27 (Supplement): 1–36.
- Botosan, C. A. 1997. Disclosure level and the cost of equity capital. *The Accounting Review* 72 (3): 323–349.
- Breiman, L., J. Friedman, R. A. Olshen, and C. J. Stone. 1984. *Classification and Regression Trees*. Boca Raton, FL: Chapman & Hall/CRC.
- . 2001. Statistical modeling: The two cultures. *Statistical Science* 16 (3): 199–215.
- Burgstahler, D., and I. Dichev. 1997. Earnings management to avoid earnings decreases and losses. *Journal of Accounting and Economics* 24: 99–126.
- Chan, L. K., N. Jegadeesh., and J. Lakonishok. 1996. Momentum strategies. *The Journal of Finance* 51 (5): 1681–1713.
- Clatworthy, M., and M. J. Jones. 2003. Financial reporting of good news and bad news: Evidence from accounting narratives. *Accounting and Business Research* 33 (3): 171–185.
- DeBondt, W. F. M., and R. Thaler. 1985. Does the stock market overreact? *The Journal of Finance* 40 (3): 793–805.
- Fama, E., and K. French. 1997. Industry costs of equity. *Journal of Financial Economics* 43: 153–193.
- Foster, G., C. Olsen, and T. Shevlin. 1984. Earnings releases, anomalies, and the behavior of security returns. *The Accounting Review* 59 (4): 574–603.
- Francis, J., K. Schipper, and L. Vincent. 2002. Expanded disclosures and the increased usefulness of earnings announcements. *The Accounting Review* 77 (3): 515–546.
- Frazier, K., R. Ingram, and M. Tennyson. 1984. A methodology for the analysis of narrative accounting disclosures. *Journal of Accounting Research* 22 (1): 318–331.
- Freedman, D. A. 1983. A note on screening regression equations. *The American Statistician* 37 (2): 152–155.
- Hand, D., H. Mannila, and P. Smyth. 2001. *Principles of Data Mining*. Cambridge, MA: MIT Press.
- Hart, R. P. 2000. *Diction 5.0, the Text-Analysis Program, Users Manual*. Thousand Oaks, CA: Scolari, Sage Publications.
- Hastie, T., R. Tibshirani, and J. Friedman. 2001. *The Elements of Statistical Learning: Data Mining, Inference and Prediction*. New York, NY: Springer-Verlag.
- Hayn, C. 1995. The information content of losses. *Journal of Accounting and Economics* 20: 125–153.
- Henry, E. 2005. Are investors influenced by the way earnings press releases are written? Working paper, Rutgers University.
- Hoskin, R. E., J. S. Hughes, and W. E. Ricks. 1986. Evidence on the incremental information content of additional firm disclosures made concurrently with earnings. *Journal of Accounting Research* 24 (Supplement): 1–32.
- Hotchkiss, E. S., and D. Strickland. 2003. Does shareholder composition matter? Evidence from the market reaction to corporate earnings announcements. *The Journal of Finance* 58 (4): 1469–1498.
- Hussainey, K., T. Schleicher, and M. Walker. 2003. Undertaking large-scale disclosure studies when AIMR-FAF ratings are not available: The case of prices leading earnings. *Accounting and Business Research* 33 (4): 275–294.
- Jegadeesh, N., and S. Titman. 1993. Returns to buying winners and selling losers: Implications for stock market efficiency. *The Journal of Finance* 48 (1): 65–91.

- Kasznik, R., and B. Lev. 1995. To warn or not to warn: Management disclosures in the face of an earnings surprise. *The Accounting Review* 70 (1): 113–134.
- Kennedy, P. 1998. *A Guide to Econometrics*. 4th edition. Cambridge, MA: MIT Press.
- Kloptchenko, A., T. Eklun, J. Karlsson, B. Back, J. Vanharanta, and A. Visa. 2004. Combining data and text mining techniques for analyzing financial reports. *Intelligent Systems in Accounting, Finance and Management* 12 (1): 29–41.
- Kothari, S. 2001. Capital markets research in accounting. *Journal of Accounting and Economics* 31 (1–3): 105–231.
- Lev, B. 1989. On the usefulness of earnings and earnings research: Lessons and directions from two decades of empirical research. *Journal of Accounting Research* 27 (3): 153–201.
- Lopez, T., and L. Rees. 2002. The effect of beating and missing analysts' forecasts on the information content of unexpected earnings. *Journal of Accounting, Auditing and Finance* 17 (2): 155–184.
- Lovell, M. C. 1983. Data mining. *The Review of Economics and Statistics* 65 (1): 1–12.
- Manning, C. D., and H. Schutze. 1999. *Foundations of Statistical Natural Language Processing*. Cambridge, MA: MIT Press.
- McCarthy, W. 1982. The REA accounting model: A generalized framework for accounting systems in a shared data environment. *The Accounting Review* 57 (3): 554–578.
- Miller, G. S. 2002. Earnings performance and discretionary disclosure. *Journal of Accounting Research* 40 (1): 173–203.
- Moskowitz, T. J., and M. Grinblatt. 1999. Do industries explain momentum? *The Journal of Finance* 54 (4): 1249–1290.
- Neter, J., M. H. Kutner, C. J. Nachsteim, and W. Wasserman. 1996. *Applied Linear Statistical Models*. 4th edition. Chicago, IL: WCB McGraw-Hill/Irwin.
- Neuendorf, K. A. 2002. *The Content Analysis Guidebook*. Thousand Oaks, CA: Sage Publications.
- Ou, J. A., and S. Penman. 1989a. Accounting measurement, price-earnings ratio, and the information content of security prices. *Journal of Accounting Research* 27 (Supplement): 111–144.
- . 1989b. Financial statement analysis and the prediction of stock returns. *Journal of Accounting and Economics* 11 (4): 295–329.
- Rencher, A. C., and F. C. Pun. 1980. Inflation of R^2 in best subset regression. *Technometrics* 22 (1): 49–53.
- Securities and Exchange Commission (SEC). 1998. *Final Rule: Plain English Disclosure*. Release No. 33-7497 and 34-39593 (File No. S7-3-97). Available at: <http://www.sec.gov/rules/final/33-7497.txt>.
- . 2003. *Final Rule: Conditions for Use of Non-GAAP Financial Measures*. Release No. 33-8176 and 34-47226 (File No. S7-43-02). Available at: <http://www.sec.gov/rules/final/33-8176.htm>.
- Shon, J. J. 2003. The relation between earnings performance and discretionary disclosure behavior in periods of bad economic news. Working paper, University of Chicago.
- Smith, M., and R. J. Taffler. 2000. The chairman's statement: A content analysis of discretionary narrative disclosures. *Accounting, Auditing & Accountability Journal* 13 (5): 624–646.