

Modeling Analysts' Recommendations via Bayesian Machine Learning

DAVID BEW, CAMPBELL R. HARVEY, ANTHONY LEDFORD,
SAM RADNOR, AND ANDREW SINCLAIR

DAVID BEW

is principal engineer at Man-AHL in Oxford, United Kingdom.
david.bew@man.com

CAMPBELL R. HARVEY

is professor of finance at Duke University in Durham, NC, and a research associate at the National Bureau of Economic Research in Cambridge, MA.
cam.harvey@duke.edu

ANTHONY LEDFORD

is chief scientist at Man-AHL and an associate at Oxford-Man Institute of Quantitative Finance in Oxford, United Kingdom.
anthony.ledford@man.com

SAM RADNOR

is senior vice president at quant PORT in London, United Kingdom.
sradnor@quantport-am.com

ANDREW SINCLAIR

is senior quantitative analyst at Realindex Investments in Sydney, Australia.
andrew.sinclair@realindex.com.au

In 2009, a unique citizen science project called the Galaxy Zoo Supernovae project was launched.¹ One of the goals of the project was to identify new supernovae (SN)—and to recruit the help of thousands of amateur astronomers. The astronomers were asked to give three levels of classification: very likely SN object, possible SN object, and not likely SN object. Determination of the true classification came from the spectrographic analysis of Caltech's Palomar Transient Factory.²

The problem arose as to how to combine the classifications. At any point in time, many astronomers may be scoring a particular object. Should we look at the average classification? Obviously the classifications are imperfect, and an average may reduce the noise. A simple majority vote (yes or no) is another possibility. However, both the majority and the average do not allow for differential skill among the classifiers. Is there a way to build a system that takes the track record of the astronomer into account? Importantly, the quality of the track record should be dynamic to allow for both improvement through time as well as fatigue.

Such a task is an ideal application of a type of machine learning called *independent*

*Bayesian classifier combination*³ (IBCC), originally defined by Ghahramani and Kim (2003). The Galaxy Zoo data were analyzed by Simpson et al. (2013) with impressive results. They found that their probabilistic model for the IBCC technique led to dramatic improvements in classification. For example, allowing for a 10% error rate, the rate of correct classification went from approximately 65% using the average to 97% using IBCC.

What does the classification of SN have to do with finance? It turns out that there are striking similarities to the problem facing an investment manager in evaluating analysts' recommendations: As in the Galaxy Zoo project, there are thousands of objects (companies) and thousands of astronomers (analysts). In both cases, the subjects do not cover all the objects (companies), but only a subset (sparsity). The classification mechanism in the Galaxy Zoo project (very likely, possible, and not likely) has an uncanny resemblance to buy, hold, or sell. In addition, it is reasonable to assume a differential degree of skill among analysts; hence, the IBCC method, given its track record in the physical and biological sciences, is a logical place to start.

The goal of our article is to apply IBCC to the I/B/E/S forecast universe to determine

¹See Lintott (2012).

²<https://www.ptf.caltech.edu/iptf>.

³Despite its name, the IBCC model does not assume independence but, instead, assumes conditional independence, which is discussed later.

whether the classifier provides information that may lead to improved investment management. We are fully aware that analysts' forecasts are a well-researched area in the academic finance and accounting literature. Indeed, Brown (2000) detailed 575 studies, many of which are focused on analysts' forecasts—and that article is 20 years out of date. A search of SSRN's Financial Economics Network and Accounting Research Network reveals over 1,000 papers dealing with analysts' forecasts.⁴

Despite the large quantity of research, ours is the first article (that we know of) to apply IBCC to the important problem of how to combine analysts' recommendations. Previous applications of IBCC in economics include work by Levenberg et al. (2013), who focused on forecasting the trend of the US nonfarm payrolls, and by Levenberg et al. (2014), who incorporated sentiment measures obtained using sentence-level language analysis. The popularity of IBCC in large-scale machine learning applications is largely due to it providing a scalable multidimensional inference procedure for combining arbitrary groups of simultaneous recommendations from multiple sources. It does this while requiring only univariate classifier learning, thereby allowing the set of sources to be easily extended. These features also make it ideal for combining analysts' forecasts.

With the potential for incorporating so many classifier sources, avoiding overfitting becomes an important consideration. Bayesian models are not as prone to overfitting as are models that require point estimates to be specified for large numbers of parameters; uncertainty about all the unknowns in a Bayesian model is described using their joint posterior probability distribution. Prediction requires integrating over this distribution, a procedure that properly accounts for diffuse knowledge about all parameters rather than requiring point values to be ascribed. The primary drawback of Bayesian models, which automatically account for parameter uncertainty, is that their use can be computationally demanding, often making them unsuitable or even impossible for real-time use. In contrast, our inference approach uses a state-of-the-art Bayesian technique called *variational approximation*, and it is extremely efficient computationally. The model we present here can be applied to learn

about each analyst individually or groups of analysts. Restrictions currently in place require that we only report at the broker level.

We realize that predicting financial outcomes remains difficult, even when expansive datasets and sophisticated machine learning models are available. Our primary aim is not to identify the best analyst or broker but to make a coherent ensemble forecast in which the weight given to each broker is driven by the length and quality of the broker's track record. In our application, the best results arise when there is agreement between broker recommendations and the forecasts obtained using IBCC. This confirmation (or reinforcement) effect, which pervades our long-only, long-short, and short-only portfolios and the various robustness analyses we perform, suggests intriguing ways for machine learning to enhance the investment processes of both quantitative and discretionary fund managers.

Our article is organized as follows. The second section discusses the data, focusing on nonstandard features such as their categorical nature, dependence structure, and sparsity (i.e., characteristics that necessitate a bespoke modeling treatment). The third section details the IBCC model and discusses important choices about priors and hyperparameters within our Bayesian framework. The fourth section explains how inference is undertaken using a state-of-the-art computationally efficient technique called variational approximation. Empirical results are presented in the fifth section, together with a range of robustness checks. Concluding remarks and some suggestions for further research are offered in a final part.

DESCRIPTION OF THE DATA MODELING PROBLEM

Our study falls within the area of machine learning known as *supervised learning*. The input data are categorical analyst recommendations about individual companies and are obtained from a large, publicly available database. Associated with each analyst recommendation is a categorical outcome variable (sometimes called a target, or *truth*, within the IBCC literature) that describes the directional price movement of the company's stock (relative to a benchmark) subsequent to the recommendation. We aim to use a modern Bayesian machine learning method to learn the relationship between these input and target data and thereby predict future price movements based on current recommendations data.

⁴Early reviews of the literature were conducted by Givoly and Lakonishok (1984), Schipper (1991), and Brown (1993). A more recent treatment was done by Bradshaw (2011).

Input Data: I/B/E/S Broker Recommendations

A vast amount of analyst data are available on both the individual stocks and the various subsectors within international equities markets. Our focus here is on recommendation data from the Thomson Reuters I/B/E/S database, a data source that covers nearly all analysts within their respective geographies and provides analyst-by-analyst recommendations for individual securities.

A *recommendation* is simply an analyst's rating for a particular company, and because different analysts use a variety of ratings schemes, each recommendation received from a contributing analyst is mapped by Thomson Reuters to one of five Standard Ratings: strong buy, buy, hold, underperform, and sell.

Several factors distinguish such data from those typically encountered in mainstream financial forecasting applications. First, unlike in standard time-series forecasting, recommendations are not observed at a fixed frequency but are event based; that is, they are observed irregularly and at largely unpredictable discrete dates. Second, instead of being quantitative forecasts on some continuous-valued scale, recommendations are categorical. This makes them better suited to a classification-based analysis than to a standard regression approach. Additionally, the recommendation database we examine has the following characteristics:

1. **Very high dimensionality:** Recommendations are received on thousands of stocks from thousands of individual analysts.
2. **Extreme sparsity:** Typically only a small number of analysts issue recommendations on any particular stock on any particular day; the rest say nothing.
3. **Dependence:** We expect analyst recommendations to be statistically dependent for a number of reasons:

A. **Cross-sectional dependence:** Contributing analysts often have exposure to correlated information sets and therefore reach the same or similar conclusions even though their decision processes are otherwise independent. This is an example of an important special case in statistics: When a multivariate random variable, $(X_1, X_2, \dots, X_m, Z)$ say, is such that $\Pr(X_1, X_2, \dots, X_m | Z) = \Pr(X_1 | Z) \times \Pr(X_2 | Z) \times \dots \times \Pr(X_m | Z)$, then the X s are said

to be conditionally independent given Z , or equivalently, the X s are independent conditional on Z . The IBCC model makes extensive use of such a conditional independence structure (see the third section).

B. **Temporal dependence:** Analyst views typically update gradually, and analysts often restate their previous recommendations. This leads to serial correlation. Group behavior among analysts can also generate serial correlation (e.g., some analysts leading opinion and others following consensus).

4. **Lack of consistency:** Although the analyst recommendations provided by I/B/E/S are recorded on the common five-category scale given earlier, for many analysts, only two of these categories are populated. Other analysts may use three of the available categories and still others all five. Although it is quite possible to deal with this inconsistency using all five categories within the IBCC model, there is little practical gain in doing so here. Thus, we group together the first two and last two Thomson Reuters Standard Ratings and relabel the original I/B/E/S analyst recommendations as buy, hold, and sell. For each I/B/E/S recommendation, we artificially label each analyst not issuing a recommendation for that stock-day pair with the category label "Missing." This means that recommendations are recorded on the following four-category scale: missing, buy, hold, and sell. Finally, we note that the distribution of buys and sells can be extremely uneven reflecting inherent biases in broker behavior.

Accounting for any one of these four characteristics within a Bayesian analysis requires detailed probabilistic modeling. Our IBCC methodology deals with all of them simultaneously and does so with a computationally rapid approach that allows the resulting system to calibrate dynamically to the prevailing environment. We also require that the prediction computations required for forecasting be feasible in real time so incoming recommendations can be responded to with minimal delay. Our Bayesian approach also allows prior beliefs to be accommodated so that the system can be guided by information from outside the observed data, should that be required.

Outcome Data: Post-Recommendation Price Movements

Unlike the input recommendations data, which are intrinsically categorical, the outcome data we seek to predict are price movements of the underlying company's stock over some future time horizon. Such price movements arise on an essentially continuous rather than categorical scale, whereas the IBCC model, which we seek to apply here, requires categorical targets. Our first step is therefore to create these categorical targets for the historical recommendation data. For consistency with the IBCC literature, these targets will be referred to as *truths*.

We first need to choose the time horizon, $\Delta\tau$, over which we are interested in predicting the movement of the stock price; for the majority of this study, we use $\Delta\tau = 60$ business days. For each analyst recommendation, we note the day it became public, s , and calculate $r_{(s,\Delta\tau)}$, which is defined as the excess return of the relevant stock over the $\Delta\tau$ period starting the next business day after s and measured relative to our benchmark return. We use this together with a relative measure of index volatility to define a categorical truth variable t for each recommendation according to

$$t = \begin{cases} 0, & \text{if } r_{(s,\Delta\tau)} \leq -5\% \times RVol_s, \\ 2, & \text{if } r_{(s,\Delta\tau)} \geq 5\% \times RVol_s, \\ 1, & \text{otherwise} \end{cases}$$

where $RVol_s$ denotes an estimator of index volatility scaled to have unit mean. Given their obvious interpretations, we refer interchangeably to the truth states $\{0, 1, 2\}$ as Price_Down, Price_Flat, and Price_Up, respectively. Clearly, the truth variable defined here has nothing to do with any broker recommendation being correct or incorrect; it is determined solely by the subsequent performance of the stock relative to the index after the recommendation. Many reasonable extensions of this truth variable definition are possible—for example, one could incorporate the market β of each underlying stock.

We restrict our attention to the period January 1, 2004, to January 1, 2013, and include only the pan-European region comprising Austria, Belgium, the Czech Republic, Cyprus, Denmark, Finland, France, Germany, Greece, Hungary, Ireland, Italy, Luxembourg,

EXHIBIT 1 The Structure of the Dataset

Stock ID	Date	Truth	Broker 1	Broker 2	...	Broker N
#1234	July 10, 2008	0	3	0		0
#5678	Feb 7, 2012	0	0	1		2
#5678	July 1, 2012	2	2	0		0
#5678	Mar 14, 2012	1	0	3		1
	:	:	:	:	:	:

Notes: Integer codes $\{0, 1, 2\}$ are used to denote the truth outcomes $\{\text{Price_Down}, \text{Price_Flat}, \text{Price_Up}\}$, respectively. The artificial recommendation label “Missing” is encoded as 0 for each noncontributing broker in each row, and the resulting recommendation set $\{\text{Missing}, \text{Hold}, \text{Sell}, \text{Buy}\}$ is encoded as $\{0, 1, 2, 3\}$, respectively. Each row contains at least one nonzero recommendation code. A very high proportion of recommendations is recorded as 0, corresponding to “Missing,” because only a small number of brokers within each row issues hold, sell, or buy recommendations.

the Netherlands, Norway, Poland, Portugal, Russia, Spain, Sweden, Switzerland, Turkey, and the United Kingdom. Our benchmark is the Dow Jones Euro Stoxx Index.⁵ Additionally, at the announcement time of each recommendation, we apply a filter to the stock universe to ensure our results are free of survivorship bias.

We group analysts by their stated corporate employer, henceforth *broker*, which gives 347 separate brokers. To be clear, the IBCC technique can be applied at the level of individual analysts or at the broker level. Because of reporting restrictions, we focus this article at the broker level.

Aggregating recommendations about the same stock that arise on the same day, we obtain the combined recommendations and truths dataset described in Exhibit 1. The dataset has 105,319 rows.⁶ If recommendations were recorded for all 347 brokers for each of the 105,319 rows there would be 36,545,693 nonzero recommendation codes, corresponding to combinations of the labels hold, sell, and buy. However, the reality is that only 116,220 of the recommendation codes in Exhibit 1 are nonzero, meaning 99.7% correspond to the label “Missing.” This demonstrates the extreme sparsity of the data object at the heart of our IBCC analysis.

⁵ Bloomberg ticker: SXXE Index.

⁶ This choice of a one-day aggregation period is arbitrary and is something we return to later. From the previous discussion about group behavior, we would expect statistical dependence between rows at this aggregation were analysts to issue recommendations on a stock prompted by others doing so.

THE IBCC MODEL: PROBABILISTIC SPECIFICATION AND CONSTRUCTION OF THE POSTERIOR

The IBCC model is a fully probabilistic model that relates a constellation of categorical inputs—in our case, the constellation of broker recommendations within each row of the data object described in Exhibit 1—and a categorical truth variable associated with those inputs.⁷

We start by specifying a probabilistic model over the categorical truth variable T . In our IBCC implementation, T takes values over states $\{0, 1, 2\}$ corresponding to Price_Down, Price_Flat, and Price_Up, respectively, and is assumed to have probability mass function

$$\Pr(T = t | \mathbf{\kappa}) = \kappa_t \text{ for } t \in \{0, 1, 2\} \quad (1)$$

where the parameter $\mathbf{\kappa} = (\kappa_0, \kappa_1, \kappa_2)$ denotes a three-vector of probabilities so that $\kappa_0 + \kappa_1 + \kappa_2 = 1$. This specification is simply saying the truths $\{0, 1, 2\}$ occur with probabilities $\mathbf{\kappa} = (\kappa_0, \kappa_1, \kappa_2)$ respectively, and that no other truth outcomes are possible. The conditioning notation in Equation 1 makes explicit that the parameter $\mathbf{\kappa}$ is assumed to be known at this stage.

The next step is to specify, for each broker, three separate distributions to describe their recommendation behavior given each possible truth. More explicitly, letting $B_k \in \{0, 1, 2, 3\}$ denote the recommendation of broker k corresponding to missing, hold, sell, and buy, respectively, for each $k \in \{1, \dots, N\}$ we require distributions for the following three conditional random variables: $B_k | T = 0$, $B_k | T = 1$, and $B_k | T = 2$. Writing T_j for the truth in row j , the IBCC model assumes, conditionally on $T_j = t$, that the B_k are independent and have probability mass functions given by

$$\Pr(B_k = b_{kj} | T_j = t, \boldsymbol{\pi}_t^{(k)}) = \pi_{t, b_{kj}}^{(k)} \text{ for } b_{kj} \in \{0, 1, 2, 3\} \quad (2)$$

where, for each truth $t \in \{0, 1, 2\}$, the parameter $\boldsymbol{\pi}_t^{(k)} = [\pi_{t,0}^{(k)}, \pi_{t,1}^{(k)}, \pi_{t,2}^{(k)}, \pi_{t,3}^{(k)}]$ denotes a four-vector of probabilities for broker k and so satisfies $\pi_{t,0}^{(k)} + \pi_{t,1}^{(k)} + \pi_{t,2}^{(k)} + \pi_{t,3}^{(k)} = 1$. This conditional specification looks complicated, but all we are doing is defining three separate four-dimensional multinomial distributions for each broker, one for each

of the possible truth outcomes. Thus, for each broker k , we have parameters $\boldsymbol{\pi}_0^{(k)}$, $\boldsymbol{\pi}_1^{(k)}$, and $\boldsymbol{\pi}_2^{(k)}$. Again, the conditioning notation in Equation 2 makes explicit that the parameters $\boldsymbol{\pi}_t^{(k)}$ are assumed known at this point.

The assumption that the broker recommendations within row j are independent conditionally on $T_j = t_j$ allows the likelihood contribution for row j to be constructed, giving

$$\Pr(T = t_j, B_1 = b_{1j}, B_2 = b_{2j}, \dots, B_N = b_{Nj}) = \kappa_{t_j} \pi_{t_j, b_{1j}}^{(1)} \pi_{t_j, b_{2j}}^{(2)} \cdots \pi_{t_j, b_{Nj}}^{(N)} = \kappa_{t_j} \prod_{l=1}^N \pi_{t_j, b_{lj}}^{(l)}$$

The IBCC model assumes all rows in the data object described in Exhibit 1 are independent, so the full likelihood, over its n distinct rows, is given by

$$\Pr(\mathbf{t}, \mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_N) = \prod_{j=1}^n \left(\kappa_{t_j} \prod_{l=1}^N \pi_{t_j, b_{lj}}^{(l)} \right) \quad (3)$$

where for notational brevity we have written $\mathbf{t} = (t_1, \dots, t_n)$ for the column of n truths in Exhibit 1 and $\mathbf{b}_k = (b_{k1}, \dots, b_{kn})$ for the recommendation column for broker k , for each $k \in \{1, \dots, N\}$.

So far we have treated the parameters of the truths and broker distributions—that is, $\mathbf{\kappa}$ and $(\boldsymbol{\pi}_0^{(k)}, \boldsymbol{\pi}_1^{(k)}, \boldsymbol{\pi}_2^{(k)})$ for $k \in \{1, \dots, N\}$, respectively—as fixed parameters. In a frequentist analysis, we would need to estimate these (e.g., by maximum likelihood and its well-established asymptotic theory) to obtain point estimates and confidence intervals. This is not the approach we adopt here. Our Bayesian analysis requires that we treat all these quantities probabilistically so that each is described according to its own prior probability distribution. Formulation of the posterior distribution then proceeds via the product of these prior distributions and the likelihood given in Equation 3, and inferences are made based on the posterior distribution alone (see Lee 2012).

Thus, we must specify priors over $\mathbf{\kappa}$ and $\boldsymbol{\pi}_0^{(k)}, \boldsymbol{\pi}_1^{(k)}$, and $\boldsymbol{\pi}_2^{(k)}$ for each $k \in \{1, \dots, N\}$. Because the truth and broker recommendation distributions are all examples of multinomial distributions, we choose to use the family of Dirichlet distributions as priors because the

⁷ Code for IBCC is available at <https://github.com/edwinrobots/pyIBCC>. This is not the code that we used for our research.

Dirichlet family is the conjugate⁸ family of priors for the multinomial distribution (for details, see Bishop 2006).

For the truth probabilities $\mathbf{\kappa} = (\kappa_0, \kappa_1, \kappa_2)$, we assume a three-dimensional Dirichlet distributed prior, that is, the continuous distribution with probability density function over domain of support $D = \{(\kappa_0, \kappa_1, \kappa_2); 0 \leq \kappa_j \leq 1, \sum_{t=0}^2 \kappa_t = 1\}$ given by $\Pr(\mathbf{\kappa} | \mathbf{v}) = C(\mathbf{v}) \prod_{t=0}^2 \kappa_t^{v_{0t}-1}$, where $C(\mathbf{v}) = \Gamma(v_{00} + v_{01} + v_{02}) / \{\Gamma(v_{00}) \times \Gamma(v_{01}) \times \Gamma(v_{02})\}$; the three-vector $\mathbf{v} = (v_{00}, v_{01}, v_{02})$ denotes a so-called *hyperparameter* (i.e., a parameter of the prior); and $\Gamma(\cdot)$ is the gamma function. Note that substituting $v_{0t} \equiv 1$ into this probability density function for each $t \in \{0, 1, 2\}$ yields a flat prior for $\mathbf{\kappa}$ over D . Similarly, for each broker recommendation B_k for $k \in \{1, \dots, N\}$ and conditional on truth $t \in \{0, 1, 2\}$ we assume $\{\pi_{t0}^{(k)}, \pi_{t1}^{(k)}, \pi_{t2}^{(k)}, \pi_{t3}^{(k)}\}$ has a four-dimensional Dirichlet prior with hyperparameters $\{\alpha_{0,t0}^{(k)}, \alpha_{0,t1}^{(k)}, \alpha_{0,t2}^{(k)}, \alpha_{0,t3}^{(k)}\}$. To condense the notation, we denote the complete set of broker recommendation probabilities conditional on each truth by $\mathbf{\Pi} = [\{\pi_{t0}^{(k)}, \pi_{t1}^{(k)}, \pi_{t2}^{(k)}, \pi_{t3}^{(k)}\}; t = 0, 1, 2; k = 1, \dots, N]$ and their corresponding hyperparameters by $\mathbf{A}_0 = [\{\alpha_{0,t0}^{(k)}, \alpha_{0,t1}^{(k)}, \alpha_{0,t2}^{(k)}, \alpha_{0,t3}^{(k)}\}; t = 0, 1, 2; k = 1, \dots, N]$.

Having now fully specified both the likelihood and the prior, we are equipped to construct the posterior distribution, which is proportional to their product, and hence satisfies

$$\Pr(\mathbf{\kappa}, \mathbf{\Pi}, \mathbf{t}, \mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_N | \mathbf{A}_0, \mathbf{v}) \propto \prod_{j=1}^n \left(\kappa_{t_j} \prod_{l=1}^N \pi_{t_j b_{jl}}^{(l)} \right) \Pr(\mathbf{\kappa} | \mathbf{v}) \Pr(\mathbf{\Pi} | \mathbf{A}_0). \quad (4)$$

This is a joint distribution in over 4,000 dimensions⁹ and incorporates information about $\mathbf{\kappa}$ and $\mathbf{\Pi}$ from both the observed data and the prior. In some IBCC applications, it is common to choose informative priors; however, we deliberately choose priors that are flat over their respective domains. This is achieved by setting the hyperparameters $v_{0t} \equiv 1$ for each $t \in \{0, 1, 2\}$, and $\alpha_{0,tj}^{(k)} \equiv 1$ for each $j \in \{0, 1, 2, 3\}$ where $\{k \in \{1, \dots, N\}; t \in \{0, 1, 2\}\}$. These flat priors ensure that only information learned from the observed truths and recommendations data,

and not our choice of priors, is driving the trading signals and allows straightforward assessment of the efficacy of our learning framework.

For the avoidance of doubt, we note that the IBCC model incorporates no sense of ordering within the category labels for either the truths or the broker recommendations. Its fundamental job is simply to learn how one set of labels (the broker recommendations) relates to the other set (the truths, which encode subsequent price outcome). Indeed, a broker that always recommends buy when the truth is Price_Down is just as informative within our IBCC implementation as a broker that always recommends sell in such cases.

The high dimensionality and data sparsity of our application mean using alternative dependence models (e.g., copulas) to capture the dependence between different brokers is computationally infeasible. The IBCC model deals with this limitation by assuming conditional independence and thereby provides a scalable and computationally efficient multidimensional inference procedure over arbitrary groups of classifiers that requires only univariate classifier learning. This key feature of the IBCC model is one of the reasons it has become popular for large-scale Bayesian machine learning applications.

VARIATIONAL BAYESIAN INFERENCE

In this section, we introduce variational Bayesian inference, an approach sometimes termed *variational Bayes*, or simply VB. See Bishop (2006, Chapter 10) and Blei, Kucukelbir, and McAuliffe (2018) for detailed treatments and Fox and Roberts (2011) for a tutorial.¹⁰ We then provide the key results of applying VB to our IBCC model. The theory is elegant, but its mathematical derivation can obscure the simplicity of the underlying approach: We approximate a multivariate distribution by a product of simpler distributions that we update iteratively to obtain the best overall approximation. In what follows, all logarithms are natural logs, that is, $\log_e(\cdot)$.

Let \mathbf{X} denote a set of observed data and \mathbf{Z} a combined set of latent (i.e., unobserved) parameters and variables. We use the generic shorthand $p(\cdot)$ to denote the probabilistic model governing whatever quantities appear inside the parentheses; for example, the joint distribution of \mathbf{X} and \mathbf{Z} is written $p(\mathbf{X}, \mathbf{Z})$. Our goal is to find a good approximation, $q(\mathbf{Z})$ say, for the posterior

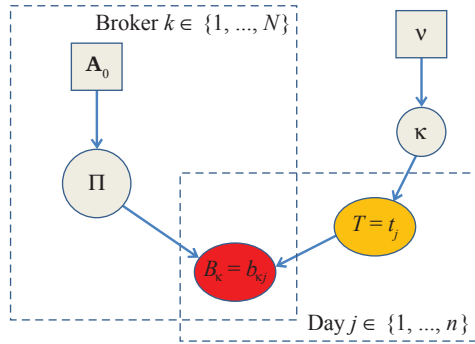
⁸ A conjugate prior is one that leads to a posterior distribution that is within the same parametric family as the prior, which therefore leads to greatly simplified Bayesian analysis. See Bishop (2006).

⁹ There are 347 brokers, each requiring three separate four-dimensional distributions, plus the three-dimensional truth distribution. In all, this makes $347 \times 4 \times 3 + 3 = 4,167$ dimensions.

¹⁰ Also, see <https://staff.aist.go.jp/bevan.jones/vb-tutorial-slides.pdf>.

EXHIBIT 2

Graphical Model of Our IBCC Implementation



Notes: Elliptical/circular nodes are variables with a distribution, whereas rectangular nodes represent hyperparameter variables that are instantiated with fixed values. The red shaded node represents recommendations, which are observed during both training and prediction. The orange shaded node represents truths, which are observed during training but have to be inferred during prediction.

$p(\mathbf{Z}|\mathbf{X})$. In our IBCC implementation, \mathbf{Z} will include the truth outcome we seek to predict (i.e., Price_Up, Price_Down, or Price_Flat; see Exhibit 2).

Noting that $q(\mathbf{Z})$ represents a probability model and therefore integrates to one, we may always write $\log p(\mathbf{X}) = \int q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z}$, where $p(\mathbf{X})$ denotes the so-called *model evidence*. Furthermore, because the definition of conditional probability gives $p(\mathbf{X}) = p(\mathbf{X}, \mathbf{Z})/p(\mathbf{Z}|\mathbf{X})$, we may substitute for $p(\mathbf{X})$ in this integral to obtain

$$\begin{aligned} \log p(\mathbf{X}) &= \int q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})} \right\} d\mathbf{Z} \\ &= \int q(\mathbf{Z}) \log \left\{ \frac{p(\mathbf{X}, \mathbf{Z})}{q(\mathbf{Z})} \times \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X})} \right\} d\mathbf{Z} \end{aligned}$$

This can be written as $\log p(\mathbf{X}) = L(q) + KL(q, p)$, where $KL(q, p) = -\int q(\mathbf{Z}) \log \{p(\mathbf{Z}|\mathbf{X})/q(\mathbf{Z})\} d\mathbf{Z}$ denotes the Kullback–Leibler¹¹ divergence (KL-divergence) between $q(\mathbf{Z})$ and $p(\mathbf{Z}|\mathbf{X})$, and $L(q) = \int q(\mathbf{Z}) \log \{p(\mathbf{X}, \mathbf{Z})/q(\mathbf{Z})\} d\mathbf{Z}$ is the negative of a quantity called the *variational free*

energy.¹² Standard properties of the KL-divergence include that it is always nonnegative and that $KL(q, p) = 0$ if and only if $q(\mathbf{Z})$ equals $p(\mathbf{Z}|\mathbf{X})$. This implies $L(q)$ is a lower bound for $\log p(\mathbf{X})$ and furthermore that this lower bound can be maximized by minimizing the KL-divergence, $KL(q, p)$, with respect to the distribution $q(\mathbf{Z})$. This is a *calculus of variations* problem.¹³

VB considers a restricted but tractable family of distributions to represent $q(\mathbf{Z})$ and then seeks the element of that family that maximizes $L(q)$. The approach we adopt involves partitioning \mathbf{Z} into m groups of variables and assuming that $q(\mathbf{Z})$ can be approximated by the factorized structure $q(\mathbf{Z}) = \prod_{i=1}^m q_i(\mathbf{Z}_i)$. This factorized version of variational approximation has its origins in physics, where it is called *mean field theory*.¹⁴ Thus, among all distributions of the form $q(\mathbf{Z}) = \prod_{i=1}^m q_i(\mathbf{Z}_i)$, we seek the distributions $q_i^*(\mathbf{Z}_i)$ that jointly maximize $L(q)$. To be clear, other than the assumed factorization structure $q(\mathbf{Z}) = \prod_{i=1}^m q_i(\mathbf{Z}_i)$, no further assumptions about $q(\mathbf{Z})$ are required.

Substituting our assumed factorization $q(\mathbf{Z}) = \prod_{i=1}^m q_i(\mathbf{Z}_i)$ into the definition of $L(q)$ given earlier and adopting the notation $q_i = q_i(\mathbf{Z}_i)$, we obtain $L(q) = \int \prod_{i=1}^m q_i \{ \log p(\mathbf{X}, \mathbf{Z}) - \sum_i \log q_i \} d\mathbf{Z}$. We now rewrite this expression to make clear how it depends on one of the individual factors, $q_j(\mathbf{Z}_j)$ say, noting that any terms not involving q_j may be treated as constant with respect to \mathbf{Z}_j . We thereby obtain

$$\begin{aligned} L(q) &= \int q_j \left\{ \int \log p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i d\mathbf{Z}_i \right\} d\mathbf{Z}_j \\ &\quad - \int q_j \log q_j d\mathbf{Z}_j + \text{Constant} \end{aligned} \quad (5)$$

We now define the new distribution $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$ by $\log \tilde{p}(\mathbf{X}, \mathbf{Z}_j) = E_{i \neq j} [\log p(\mathbf{X}, \mathbf{Z})] + c$, where c is a normalization constant and $E_{i \neq j}[\cdot]$ denotes expectation with respect to all q_i distributions for $i \neq j$ so that

¹² To avoid the possibility of misinterpretation, for clarity we remark that $L(q)$ is not the likelihood function. Writing $-L(q) = -\int q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}) d\mathbf{Z} - \int q(\mathbf{Z}) \log \{1/q(\mathbf{Z})\} d\mathbf{Z}$, we obtain an energy term minus an entropy term, which is why it is called the free energy. See Sato (2001).

¹³ Standard calculus allows functions to be optimized, where a function is a map that takes the value of some variable as input and returns the value of the function as output. Calculus of variations allows functionals to be optimized rather than functions, where functionals are maps that take functions as inputs.

¹⁴ See Parisi (1988).

$E_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})] = \int \log p(\mathbf{X}, \mathbf{Z}) \prod_{i \neq j} q_i d\mathbf{Z}_i$. Careful inspection of Equation 5 now shows that $L(q)$ is simply the negative KL-divergence between $q_j(\mathbf{Z}_j)$ and $\tilde{p}(\mathbf{X}, \mathbf{Z}_j)$, which is minimized by taking $q_j(\mathbf{Z}_j) = \tilde{p}(\mathbf{X}, \mathbf{Z}_j)$. Thus, keeping q_i constant for each $i \neq j$, we have that maximizing $L(q)$ over all possible distributions $q_j(\mathbf{Z}_j)$ is achieved by taking $\log q_j^*(\mathbf{Z}_j) = E_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})] + c$, where c denotes a normalizing constant. This key result provides the basis for application of variational methods.

The set of equations $\log q_j^*(\mathbf{Z}_j) = E_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})] + c$ for each $j \in \{1, \dots, m\}$ provides conditions for the maximum of $L(q)$ subject to the assumed factorization $q(\mathbf{Z}) = \prod_{i=1}^m q_i(\mathbf{Z}_i)$. However, these equations do not provide an explicit solution because the expression for each $q_j^*(\mathbf{Z}_j)$ involves taking expectation with respect to the other $q_i^*(\mathbf{Z}_i)$ distributions for $i \neq j$. To solve these equations, we proceed iteratively. First, each $q_i(\mathbf{Z}_i)$ distribution is initiated—for example, with parameters chosen broadly to match moments of the observed data. Then we cycle through each $j \in \{1, \dots, m\}$, updating $q_j(\mathbf{Z}_j)$ by evaluating $E_{i \neq j}[\log p(\mathbf{X}, \mathbf{Z})]$ using the current estimates of $q_i(\mathbf{Z}_i)$ for each $i \neq j$. Convergence to a local maximum is guaranteed because of certain convexity properties of $L(q)$ with respect to the factors $q_i(\mathbf{Z}_i)$ (see Boyd and Vendenbergh 2004). Furthermore, in the particular case of our IBCC implementation, because all the factors we have chosen are of the exponential family type (Bernardo and Smith 1994), this maximum can be shown to be the global maximum within the family of factorized distributions.

Variational Inference for Our IBCC Implementation

Our IBCC application deviates from those of Kim and Ghahramani (2012) and Simpson et al. (2013) in three key ways. First, we intend to perform online forecasting, so temporal consistency requires running the model using only information that is already available at the time of each forecast. Second, with the exception of truths corresponding to recommendations made within the most recent $\Delta\tau$ period, all truths within our training data are completely observed because they are based on publicly available price data. In contrast, for the Galaxy Zoo project, the truth data were largely missing. Finally, our primary interest is the predictive distribution $\Pr(T = t | B_1 = b_1, B_2 = b_2, \dots, B_N = b_N)$, rather than the posterior, because we wish to forecast the truth

outcome conditional on, for example, today's constellation of broker recommendations.¹⁵

Although it is possible to extend the IBCC model to include explicit temporal structure, as done by Simpson et al. (2013), our approach is based on calibrating their simpler static model to a dataset that updates as time evolves. Specifically, we truncate the observed data, comprising the time-stamped recommendations and truths, at a sequence of evaluation dates, ensuring additionally that a buffer of duration $\Delta\tau$ is incorporated between the last admitted training observations and the onset of prediction. For each training data set so created, we seek to calculate the predictive distribution $\Pr(T = t | B_1 = b_1, B_2 = b_2, \dots, B_N = b_N)$ for each constellation of broker recommendations that arises until the next evaluation date. Learning remains halted over this prediction phase, so each constellation of analyst recommendations we use in prediction is treated individually. All our findings are obtained using this rolling out-of-sample scheme.

For each evaluation date, we undertake both expanding-window and moving-window analyses. The expanding-window analysis admits all data from January 1, 2004, up to the evaluation date, whereas the moving-window analysis admits only data within a three-year lookback from each evaluation date. In principle, the evaluation dates could be chosen to index each business day; however, for practical reasons¹⁶ we set them quarterly, to the first day of March, June, September, and December.

Let index $i \in (1, \dots, n_i)$ denote the rows of the training data, renumbered as required for the rolling window case. Because all the recommendations and truths are observed for these training data and because we chose conjugate Dirichlet priors for both $\boldsymbol{\kappa}$ and $\boldsymbol{\Pi}$, standard properties of the multinomial-Dirichlet family (see Bishop 2006) give the following:

1. The posterior of $\boldsymbol{\kappa}$ is a Dirichlet distribution with parameter $\mathbf{v}^* = (v_0^*, v_1^*, v_2^*)$, where $v_j^* = v_{0j} + N_j$

¹⁵ The predictive distribution we seek, sometimes called the *posterior predictive distribution*, is defined by the multivariate integral $\int \Pr(T = t | B_1 = b_1, B_2 = b_2, \dots, B_N = b_N, \boldsymbol{\kappa}, \boldsymbol{\Pi}) \Pr(\boldsymbol{\kappa}, \boldsymbol{\Pi}) d\boldsymbol{\kappa} d\boldsymbol{\Pi}$, where $\Pr(\boldsymbol{\kappa}, \boldsymbol{\Pi})$ denotes the posterior distribution of $(\boldsymbol{\kappa}, \boldsymbol{\Pi})$, which depends implicitly on the training data.

¹⁶ Risk managers tend to have a preference for models with parameters that remain static for reasonable periods rather than models in which parameters change on a daily basis.

and N_j denotes the number of occurrences of truth j in the training data for $j \in \{0, 1, 2\}$. The $\mathbf{v}_j^* = \mathbf{v}_{0,j} + N_j$ formula is often referred to as the *prior counts plus data counts* updating relationship for the multinomial-Dirichlet family.

2. The posterior of $\pi_t^{(l)}$ is a Dirichlet distribution with parameters $(\alpha_{t0}^{(l*)}, \alpha_{t1}^{(l*)}, \alpha_{t2}^{(l*)}, \alpha_{t3}^{(l*)})$, where $\alpha_{tb}^{(l*)} = N_{tb}^{(l)} + \alpha_{0,tb}^{(l)}$, and $N_{tb}^{(l)}$ denotes the number of recommendations of type $b \in \{0, 1, 2, 3\}$ made in the training data by broker $l \in \{1, \dots, N\}$ for each truth $t \in \{0, 1, 2\}$. We let \mathbf{A}^* denote the collection of all these posterior parameters.

Our procedure for approximating the predictive distribution $\Pr(T=t|B_1=b_1, B_2=b_2, \dots, B_N=b_N)$ starts by considering $\Pr(\mathbf{\kappa}, \mathbf{\Pi}, t, b_1, b_2, \dots, b_N | \mathbf{A}^*, \mathbf{v}^*)$. This has the same structure as the individual data terms in Equation 4 except that now the truth t is unobserved, and $(\mathbf{A}^*, \mathbf{v}^*)$ denotes the ensemble of posterior parameters given earlier. Thus, $\log \Pr(\mathbf{\kappa}, \mathbf{\Pi}, t, b_1, b_2, \dots, b_N | \mathbf{A}^*, \mathbf{v}^*)$ is of the form

$$\sum_{j=0}^2 I(t=j) \left(\log \kappa_j + \sum_{l=1}^N \log \pi_{jb_l}^{(l)} \right) + \log \Pr(\mathbf{\kappa} | \mathbf{v}^*) + \log \Pr(\mathbf{\Pi} | \mathbf{A}^*) + \text{Constant}. \quad (6)$$

Here, we have introduced the indicator function $I(\cdot)$, defined by $I(t=j) = 1$ if $t=j$ and $I(t=j) = 0$ otherwise, because it will be convenient later. To reduce clutter, we drop the dependence on $(b_1, \dots, b_N, \mathbf{A}^*, \mathbf{v}^*)$ from our notation. We therefore represent the latent variables and parameters by $\mathbf{Z} = (t, \mathbf{\kappa}, \mathbf{\Pi})$. We assume $q(\mathbf{Z})$ factorizes as $q(t, \mathbf{\kappa}, \mathbf{\Pi}) = q(t)q(\mathbf{\kappa}, \mathbf{\Pi})$. This is the only assumption we need to make; several further simplifications arise because of the structure of the IBCC model. For example, Equation 6 shows that the terms involving $\mathbf{\kappa}$ and $\mathbf{\Pi}$ can be separated, which implies the additional factorization $q^*(\mathbf{\kappa}, \mathbf{\Pi}) = q^*(\mathbf{\kappa})q^*(\mathbf{\Pi})$.

We start by initializing the distributions for $\mathbf{\kappa}$ and $\pi_t^{(l)}$ with their posterior distributions, that is, the Dirichlet distributions with parameters \mathbf{v}^* and \mathbf{A}^* given earlier. To obtain $q^*(t)$, we need to evaluate $\log q^*(t) = E_{\mathbf{\kappa}, \mathbf{\Pi}}[\log p(t, \mathbf{\kappa}, \mathbf{\Pi})] + \text{Constant}$. Extracting the relevant terms from Equation 6, we obtain $\log q^*(t) = E_{\mathbf{\kappa}} \log \kappa_t + \sum_{l=1}^N E_{\pi_t^{(l)}} \log \pi_{tb_l}^{(l)} + \text{Constant}$. Standard properties of the Dirichlet distribution (e.g., Bishop 2006) give $E_{\mathbf{\kappa}} \log \kappa_t = \Psi(\mathbf{v}_t^*) - \Psi(\sum_{j=0}^2 \mathbf{v}_j^*)$ and

$E_{\pi_t^{(l)}} \log \pi_{tb_l}^{(l)} = \Psi(\alpha_{tb_l}^{(l*)}) - \Psi(\sum_{s=0}^3 \alpha_{ts}^{(l*)})$, where $\Psi(\cdot)$ denotes the DiGamma function.¹⁷ Next, defining the terms $\log \rho_t = \Psi(\mathbf{v}_t^*) - \Psi(\sum_{j=0}^2 \mathbf{v}_j^*) + \sum_{l=1}^N [\Psi(\alpha_{tb_l}^{(l*)}) - \Psi(\sum_{s=0}^3 \alpha_{ts}^{(l*)})]$, where b_1, \dots, b_N denote the observed broker recommendations for the prediction, we therefore obtain $q^*(t) = \rho_t / (\rho_0 + \rho_1 + \rho_2)$. This expression for $q^*(t)$ provides our initial estimate of $\Pr(T=t|B_1=b_1, B_2=b_2, \dots, B_N=b_N)$.

Deriving $q^*(\mathbf{\kappa})$ and $q^*(\mathbf{\Pi})$ requires taking expectations with respect to this newly calculated $q^*(t)$ distribution. We start by extracting the terms involving $\mathbf{\kappa}$ from Equation 6. Recalling that for any event X , the expectation of $I(X)$ is $\Pr(X)$, we obtain $\log q^*(\mathbf{\kappa}) = \sum_{j=0}^2 q(t=j) \log \kappa_j + \sum_{j=0}^2 (\mathbf{v}_j^* - 1) \log \kappa_j + \text{Constant}$. Gathering together the $\log \kappa_j$ terms in this expression shows $q^*(\mathbf{\kappa})$ to be Dirichlet distributed with parameters $\mathbf{v}_j = \mathbf{v}_j^* + q(t=j)$ for $j \in \{0, 1, 2\}$. This formula for iterating the $\mathbf{\kappa}$ distribution is similar in structure to the prior counts plus data counts relation noted previously, except that now the prior over each forecasting period is the posterior obtained at the relevant evaluation date, and the counts for each truth class are replaced with their expected values; that is, $E_t I(t=j) = q(t=j)$ for $j \in \{0, 1, 2\}$.

We essentially repeat this argument to obtain the update equations for $q^*(\mathbf{\Pi})$. First, because the $\pi_j^{(l)}$ terms in Equation 6 are separate for each truth $j \in \{0, 1, 2\}$ and each broker $l \in \{1, \dots, N\}$, we obtain the further factorization $q^*(\mathbf{\Pi}) = \prod_{l=1}^N \prod_{j=0}^2 q^*(\pi_j^{(l)})$. Extracting the $\pi_j^{(l)}$ terms and taking expectation with respect to $q^*(t)$ thereby yields $\log q^*(\pi_j^{(l)}) = \sum_{t=0}^2 q^*(t=j) \sum_{l=1}^N \log \pi_{tb_l}^{(l)} + \sum_{s=0}^3 (\alpha_{js}^{(l*)} - 1) \log \pi_{js}^{(l)} + \text{Constant}$. Gathering together terms in $\log \pi_{jb}^{(l)}$ now shows $q^*(\pi_j^{(l)})$ to be Dirichlet distributed with parameters $\alpha_{jb}^{(l)} = q(t=j)I(b=b_l) + \alpha_{jb}^{(l*)}$ for $b \in \{0, 1, 2, 3\}$ and $l \in \{1, \dots, N\}$. As before, these equations for iterating the $\mathbf{\Pi}$ distributions have the same prior counts plus expected counts interpretation.

Having updated both $q^*(\mathbf{\kappa})$ and $q^*(\mathbf{\Pi})$, we now use these distributions to obtain the next update of $q^*(t)$, and the whole scheme is iterated until convergence is obtained. The truth distribution that results is the VB

¹⁷ If the d -dimensional variable $\mathbf{X} = (X_1, \dots, X_d)$ is Dirichlet distributed with parameter (μ_1, \dots, μ_d) , then $E(\log X_i) = \Psi(\mu_i) - \Psi(\sum_{j=1}^d \mu_j)$ for each $i \in \{1, \dots, d\}$, where $\Psi(\cdot)$ denotes the DiGamma function, which is defined as $\Psi(z) = \frac{d}{dz} \log \Gamma(z)$, where $\Gamma(\cdot)$ denotes the gamma function.

approximation to $\Pr(T = t | B_1 = b_1, B_2 = b_2, \dots, B_N = b_N)$. In practice, convergence is achieved rapidly.

Although we have expressed the method in terms of a single prediction, in practice the calculations can be undertaken in parallel, allowing efficient prediction of the truth distribution for multiple constellations of broker recommendations. We remark that although the VB iteration scheme is operationally similar to the update procedure of the expectation-maximization (EM) algorithm,¹⁸ the VB and EM algorithms do very different things: EM obtains the maximum likelihood (i.e., point) estimate of a parameter, whereas VB provides a global approximation of the distribution.

From Predictive Probabilities to Decisions

The outputs of the previous procedure are the estimated truth probabilities, (q_0, q_1, q_2) say, for Price_Down, Price_Flat, and Price_Up, respectively, for each out-of-sample constellation of broker recommendations. Even when these predictive probabilities have been calculated, one still requires a decision rule—that is, a rule to decide what, if any, action to take.

We restrict our attention to the discrete set of actions Go_Short, No_Trade, and Go_Long.¹⁹ It is tempting to choose one of these actions according to whichever of Price_Down, Price_Flat, or Price_Up has the highest predictive probability (HPP). Unfortunately, this HPP rule, which chooses Go_Short if $q_0 > \max(q_1, q_2)$, Go_Long if $q_2 > \max(q_0, q_1)$, and No_Trade otherwise, is not selective enough and results in too many Go_Long actions. This behavior is unsurprising because the underlying training dataset contains unadjusted biases; analysts typically issue more buy recommendations than hold or sell, and there are more Price_Up labels than Price_Down or Price_Flat.²⁰

Recalling that q_t is an estimate of the conditional probability $\Pr(T = t | B_1 = b_1, B_2 = b_2, \dots, B_N = b_N)$, our preferred decision rule is to take the HPP action only

when q_t exceeds the current estimate of the unconditional probability of $T = t$, which is κ_t . This simple extension of the HPP rule ensures a Go_Long (Go_Short) decision arises only when knowledge of the observed constellation of broker recommendations $B_1 = b_1, B_2 = b_2, \dots, B_N = b_N$ boosts the estimated probability of Price_Up (Price_Down) relative to the background level observed within the training data.

Our default decision rule is the $c = k = 1$ case of the more general decision rule summarized as follows:

Decision	Trigger Condition
Go_Short	$q_0/\kappa_0 > c$ and $q_0 > k \max(q_1, q_2)$
No_Trade	otherwise
Go_Long	$q_2/\kappa_2 > c$ and $q_2 > k \max(q_0, q_1)$

Both parameters, c and k , affect the selectivity of this trading rule, but their effects are different and somewhat complementary. The parameter c relates to comparison of the conditional and unconditional probabilities of each truth outcome. Thus, increasing c while keeping $k = 1$ fixed means the value of the information imparted by the broker recommendations needs to be higher for a Go_Long (Go_Short) decision to arise. In contrast, the condition involving parameter k relates to the relationship among the three conditional truth probabilities, q_0, q_1 , and q_2 , but does not involve the unconditional probabilities. Thus, increasing k while keeping $c = 1$ fixed raises the threshold required for HPP decision making to produce a Go_Long (Go_Short) outcome; simply being the largest value of q_0, q_1 , and q_2 is no longer sufficient.

EMPIRICAL RESULTS AND ROBUSTNESS CHECKS

The results are based on grouping the analysts by broker (i.e., their stated corporate employers or affiliation). Learning is undertaken at this broker level and is achieved by integrating information over all the stocks and all the analysts affiliated with that broker. It is possible to implement IBCC on different types of groupings—or even by individual analysts. Such information pooling is a powerful feature of the IBCC model and Bayesian approaches more generally (e.g., providing protection against overfitting). Finer aggregations than this are possible; for example, learning could be

¹⁸ See Dempster, Laird, and Rubin (1977) and Tanner (1996).

¹⁹ Many alternatives to our discrete choice rule are possible here. For example, the calculated (q_0, q_1, q_2) probabilities could be used to derive weights on a continuous long-short scale.

²⁰ In the Galaxy Zoo project, Simpson et al. (2013) subsampled to adjust for class imbalance. We chose not to do this, instead developing a model that reflects the probabilistic structure of the observed dataset, including its biases, and dealing with these biases using an extension of the HPP decision rule.

undertaken at the Global Industry Classification Standard sector or subsector level within each broker or even at the individual analyst level, where sufficiently detailed tracking information exists to follow an analyst's career between brokers. There is, of course, a complexity penalty for finer aggregations—more model components to infer based on the same amount of data. We do not report on such aggregations here.

Another feature of our IBCC implementation is its ability to combine multiple simultaneous recommendations for each stock without the need for extra parameters. To exploit this, in the backtest simulations reported later, recommendations are aggregated over a lookback of 30 calendar days, a process that increases the number of concurrent recommendations within the rows of the training data. This procedure is best understood by considering a single stock: When a new recommendation appears, we simply look back and find the latest recommendations from the other brokers within a 30-day window and group them together in a single row of the data. Further examinations (not reported) show the impact of this choice of lookback window to be minimal.

Our standard approach is to estimate the IBCC model on a three-year period of in-sample data and then apply it out of sample to the recommendations that arise over the subsequent quarter. We then either expand or roll forward the in-sample period to include the next quarter, always applying the new fit out of sample to the following unused quarter of data. The default decision rule we use is the $c = k = 1$ case of the rule given previously. The impact of varying the parameters c and k is examined later.

We benchmark IBCC performance against a scheme that does no learning but simply aims to follow each broker's recommendations. This broker-following benchmark is referred to as *Brok_Flw* in the exhibits that follow and allows assessment of the value added by IBCC.

The *Brok_Flw* benchmark is constructed as follows:

1. For every buy recommendation, we create a signal of +1 that lasts from the day following the recommendation for 60 business days.
2. Likewise, for every sell recommendation we create a signal of -1.

3. These signals are summed within a stock, both across the multiple brokers and across multiple recommendations from the same broker.
4. The resulting signal is capped/floored at ± 10 .
5. For long-only portfolios, only underlying long recommendations are included, and conversely for short-only portfolios.
6. Each portfolio's positions are rebalanced on a daily basis to maintain a gross exposure of \$100; that is, $position_{it} = signal_{it} / \sum_i |signal_{it}|$, where the sum in this normalization is across all contemporaneous positions, both long and short.

The following nomenclature is used in presenting the results:

- **Brok_Flw_LS**: This is the broker-following benchmark described previously. We ignore recommendations in which there are simultaneous buys and sells for the same stock from different brokers.
- **IBCC_Rol_LS**: Here we apply the IBCC algorithm, fitting on a three-year rolling window, with both long and short positions.
- **IBCC_Exp_LS**: As noted earlier, but now the estimation is performed on an expanding window.
- **Both_Rol_LS**: *Both* here denotes that we only take a position if the IBCC recommendation and the raw *Brok_Flw* signal agree at the individual broker level. This prevents IBCC from reversing broker recommendations. Estimation is performed on a three-year rolling window.
- **Both_Exp_LS**: As noted earlier, but using the IBCC model on an expanding window.

Here, L (S) is used in place of LS when only long (short) positions are allowed. In all cases, the gross exposure is normalized to \$100.²¹ This means that net exposure for the LS portfolio is time varying according to the relative number of long and short recommendations. In particular, the LS results in the exhibits cannot be imputed from the separate L and S short results.

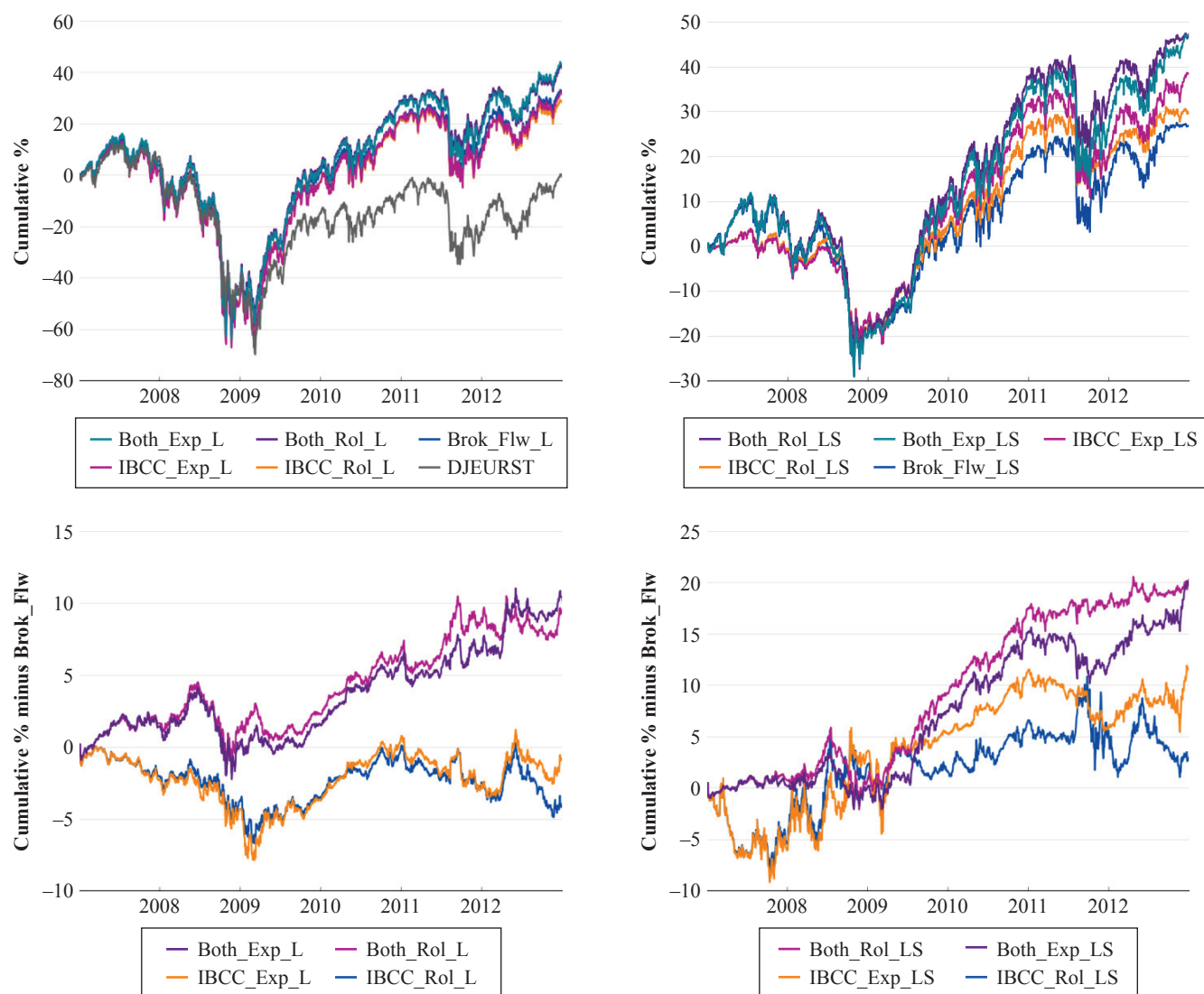
The reference index used for the intercept and slope estimates, α and β , reported in the following is the Euro Stoxx,²² the same index we used in defining

²¹ Gross exposure is defined as $\sum |pos_i|$, where pos_i is the position in the i^{th} market, in US dollars.

²² DJEURST in Thomson Reuters notation.

EXHIBIT 3

Performance of Long-Only Models (left) and Long-Short Models (right)



Notes: Outright performance is shown in the top panels, whereas the bottom panels show performance relative to the relevant *Brok_Flw** benchmark. Note the vertical axes do not share a common scale.

the truths for the training data. This is based on a liquid subset of around 300 Eurozone stocks from the STOXX Europe 600. This index had an average return close to zero over the 2007–2012 period, so the reported alphas are similar to the outright returns. Returns on short portfolios are reported assuming that all stocks have been borrowed and sold short; however, transaction and borrowing costs are not included in the results.

Exhibit 3 shows the performance of the long-only and long-short portfolios, both in terms of their outright performance and their performance relative to the relevant *Brok_Flw** benchmark. All long portfolios struggled during the global financial crisis (GFC) but comfortably outperformed the DJEURST index from 2009 onward. The long IBCC strategies remain broadly in line with the *Brok_Flw_L* benchmark, with best long performance arising for the strategies labeled *Both*.

EXHIBIT 4

Performance Statistics for the Various Long-Only, Short-Only, and Long-Short Models for the Period 2007–2012

Side	Model	Mean	Vol	Alpha	Alpha <i>t</i> -Stat	Beta	Beta <i>t</i> -Stat	Turnover
Long-Only	Brok_Flw_L	5.43	24.18	5.47	2.73	1.01	26.97	5.75
	IBCC_Rol_L	4.77	24.66	4.91	2.09	1.02	23.36	5.74
	IBCC_Exp_L	5.30	24.89	5.39	2.27	1.03	23.63	5.68
	Both_Rol_L	6.99	24.51	7.06	2.75	1.00	20.18	6.13
	Both_Exp_L	7.13	24.71	7.28	2.84	1.01	20.47	6.07
Short-Only	Brok_Flw_S	−0.51	24.96	−0.13	−0.05	−1.03	−29.42	6.38
	IBCC_Rol_S	−3.38	25.24	−3.23	−1.46	−1.05	−34.83	6.15
	IBCC_Exp_S	−3.54	24.85	−3.53	−1.60	−1.03	−34.76	6.18
	Both_Rol_S	2.99	25.98	3.45	1.06	−1.03	−21.81	7.12
	Both_Exp_S	2.11	25.71	2.46	0.79	−1.03	−25.39	7.04
Long-Short	Brok_Flw_LS	4.54	13.92	4.65	2.09	0.52	10.88	6.50
	IBCC_Rol_LS	5.07	11.01	5.30	2.69	0.39	10.63	7.23
	IBCC_Exp_LS	6.50	12.66	6.64	3.35	0.47	13.46	7.06
	Both_Rol_LS	7.99	15.43	8.15	3.18	0.56	11.11	6.32
	Both_Exp_LS	7.88	16.00	8.09	3.12	0.59	11.69	6.29

Notes: The reference index used for the α and β calculations is the Euro Stoxx, the same index used for defining the truths in the training data. The alpha values are annualized. Turnover denotes a measure of the volume traded by each portfolio on a standardized scale that allows meaningful comparison between portfolios.

For the long–short portfolios, there is no corresponding LS index, but all IBCC portfolios outperform the Brok_Flw_LS benchmark. Again, the portfolios labeled Both provide the strongest performance. Investing only when both the IBCC model and the underlying broker recommendations agree suggests a straightforward and intriguing way this machine learning application may assist investment management. No consistent benefit of fitting with rolling or expanding data windows is observed in these results.

Results for all the long-only, long–short, and short-only portfolios are tabulated in Exhibit 4, and a yearly breakdown is provided in Exhibit 5. The Brok_Flw_S benchmark and both of the short IBCC strategies are loss making, so we do not focus on their outright performance. The more interesting point is that the short portfolios labeled Both again perform better, repeating the outperformance pattern seen earlier in the long-only and long–short cases. The relative performance chart for the short-only portfolios is given in Exhibit 6 and shows the outperformance of the Both portfolios to be reasonably consistent over the post-GFC period.

Robustness Checking—Impact of Firm Liquidity

A potentially serious concern is that our IBCC procedure might be favoring recommendations from brokers who recommend smaller, less well-known stocks and thus may be inadvertently accessing a size bias. A quick check of Exhibit 7, for example, shows that Brok_Flw_L holds more stocks over \$25 billion than does IBCC.

In an attempt to control for this effect, we split the stock universe in half by market capitalization. We rank the original universe of liquid stocks by market capitalization and form a large-half backtest by including only the largest half of these stocks; in the small-half backtest, we only include the smallest half. This determination is made each month and is implemented with a five-business-day lag in an effort to reduce short-term timing effects. In the subsequent backtesting, we use these reduced universes both for the fitting of the IBCC models and subsequently for their assessment on the usual rolling out-of-sample basis. The overall number of recommendations in the two backtests is shown in Exhibit 8. The split is surprisingly even.

EXHIBIT 5

Calendar Year Performance for Long-Only, Short-Only, and Long-Short Portfolios from 2007 to 2012 Inclusive (expressed as percentage)

Side	Year	Brok_Flw	IBCC_Rol	IBCC_Exp	Both_Rol	Both_Exp	Euro Stoxx
Long-Only	2007	4.37	2.37	2.11	6.07	5.83	7.51
	2008	-47.35	-49.06	-49.38	-47.82	-48.56	-51.09
	2009	44.20	44.45	44.80	45.25	45.47	28.84
	2010	20.29	23.92	24.79	25.00	25.04	5.84
	2011	-9.99	-12.79	-13.56	-8.00	-9.33	-11.91
	2012	20.69	19.42	22.65	20.94	23.81	20.63
Short-Only	2007	5.39	-0.32	-0.99	7.07	8.83	7.51
	2008	46.59	41.07	40.97	41.64	41.26	-51.09
	2009	-42.88	-46.24	-43.83	-39.54	-38.96	28.84
	2010	-12.20	-13.29	-12.73	-7.41	-5.05	5.84
	2011	20.76	18.86	18.04	33.24	25.15	-11.91
	2012	-20.70	-20.14	-22.42	-17.26	-18.74	20.63
Long-Short	2007	5.04	-0.36	-0.99	5.87	5.29	7.51
	2008	-24.52	-16.82	-14.78	-24.42	-24.59	-51.09
	2009	22.66	21.94	24.22	30.99	29.98	28.84
	2010	16.58	21.62	22.81	24.83	24.33	5.84
	2011	-4.83	-6.51	-10.37	-4.03	-7.83	-11.91
	2012	12.00	10.18	17.62	14.12	19.51	20.63

Notes: The figures quoted are the sum of each year's daily returns. For reference, Euro Stoxx returns are shown in the right-hand column.

Backtest performance is shown in Exhibit 9 for the long-only and short-only cases,²³ and the distributions of market capitalization for the two subportfolios are shown in Exhibit 10. We conclude that IBCC is able to add value to plain I/B/E/S estimates in both large- and small-capitalization subportfolios and that the efficacy of the algorithm is not driven by a size bias.

Robustness Checking—Selectivity of the Trading Rule

We examine the impact of changing the selectivity of the trading rule so that only recommendations with progressively higher levels of conviction produce trades. The IBCC procedure remains identical to that used before; the only changes are to the values of the parameters c and k within the decision rule. This also provides a principled way to control the number of open positions. Recall that c may be

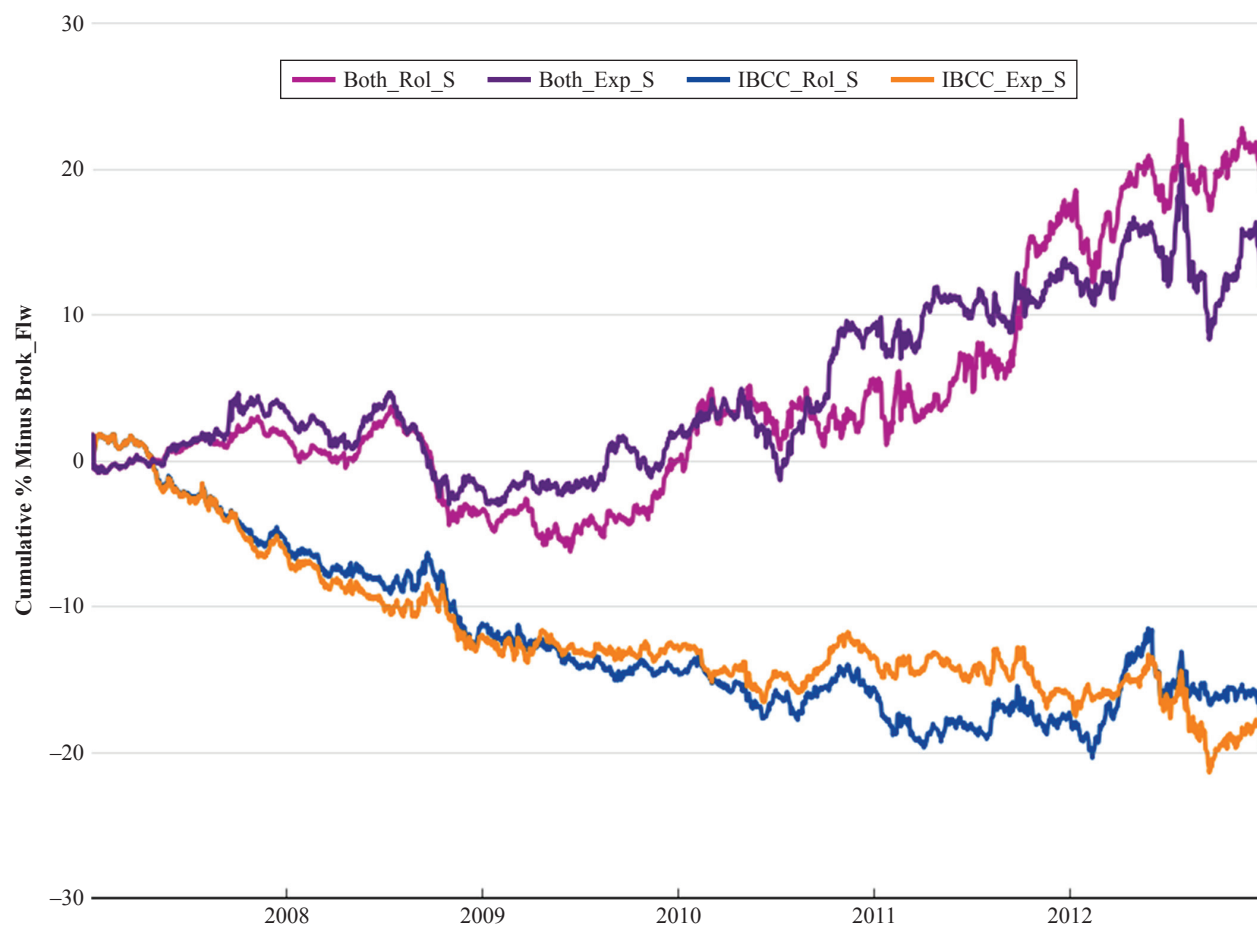
interpreted as a threshold on the information content needed within the observed constellation of broker recommendations to generate a trade. In contrast, $k > 1$ raises the threshold required for HPP decision making to produce a Go_Long (Go_Short) outcome; simply being the largest value of q_0 , q_1 , and q_2 is no longer sufficient.

We examine the impact of varying c and k separately; for space considerations, we examine only the long-only portfolios. Exhibit 11 shows the results of varying c while holding $k = 1$, and Exhibit 12 shows the results of varying k while holding $c = 1$. The results for varying c while holding $k = 1$ suggest some strengthening of both the alpha and beta as c is raised, in particular for the Both_Exp_L results. The results for varying k while holding $c = 1$ show a milder effect. As might be expected, we observed that the turnover increases as the decision rules become more selective, although the effect is mild compared to the baseline $c = k = 1$ case.

²³ We do not expect bottom-up broker recommendations to yield effective market-timing portfolios, so we do not explore the long-short case here for brevity.

EXHIBIT 6

Performance of the Short-Only Models Relative to the Brok_Flw_S Benchmark



Note: The portfolios labeled Both provide the best performance, as was the case for the long-only and long-short portfolios.

EXHIBIT 7

Size Tilts for the Different Portfolios

Model	Mega Cap >\$25 Billion	Large Cap \$10 Billion to \$25 Billion	Mid Cap \$2 Billion to \$10 Billion	Small Cap \$250 Million to \$2 Billion	Micro Cap <\$250 Million	Missing Data
Brok_Flw_L	23.1	20.8	48.7	7.2	0.0	0.0
IBCC_Rol_L	17.2	19.9	53.0	9.9	0.0	0.0
IBCC_Exp_L	16.3	19.8	53.8	9.9	0.0	0.0
Both_Rol_L	18.3	19.3	53.2	9.2	0.0	0.0
Both_Exp_L	17.5	19.3	54.0	9.2	0.0	0.0
Brok_Flw_S	15.9	20.6	52.5	10.9	0.0	0.1
IBCC_Rol_S	23.5	19.7	47.7	8.8	0.0	0.0
IBCC_Exp_S	22.8	19.6	48.7	8.6	0.1	0.1
Both_Rol_S	14.8	19.6	53.4	11.7	0.0	0.1
Both_Exp_S	14.2	19.9	54.2	11.2	0.1	0.1

Note: Exhibit shows the sum of absolute positions by market capitalization bucket, averaged across time.

EXHIBIT 8

Number of Recommendations after Bisecting the Universe of Stocks by Market Capitalization

Universe	Number of Recommendations	As a Percentage
Large Half	58,466	56%
Small Half	45,316	44%

Robustness Checking—Sensitivity to Truth Threshold

A threshold of 5% was used in the truth definition given by

$$t = \begin{cases} 0, & \text{if } r_{(s,\Delta t)} \leq -5\% \times RVol_s, \\ 2, & \text{if } r_{(s,\Delta t)} \geq 5\% \times RVol_s, \\ 1, & \text{otherwise} \end{cases}$$

This value has been used throughout. Here we explore varying this parameter between 1% and 10%, keeping everything else the same. Results are summarized in Exhibit 13 for the long-only and short-only portfolios, where for brevity we quote results only for the Both portfolios. Unreported results show that IBCC_* consistently underperforms Brok_Flw and Both_*, consistent with our previous findings.

We find that

- As before, the Both portfolios outperform the relevant Brok_Flw_* benchmark at all threshold settings for both long-only and short-only.
- There seems to be a sweet spot for thresholds within the 4%–6% range for the long-only portfolios, particularly in terms of the t -statistic for alpha, which broadly measures the consistency of the outperformance.
- In the case of short-only portfolios, a tighter threshold of around 2%–3% gives slightly better results, although nothing obtains statistical significance. One possible explanation is that the smaller number of short recommendations leads to greater sampling error in assessing a broker's short efficacy, and allocating more Price_Down truths may mitigate this.

Robustness Checking—Sensitivity to Holding Period

Here we explore the sensitivity to the arbitrary 60-day holding period that has been used throughout. For brevity, we quote results for just the long-only portfolios.

From Exhibit 14 it is reasonably clear that

- Shorter holding periods give stronger performance.
- Shorter holding periods increase turnover.
- The Both portfolios again are the strongest performers for all horizons.
- Pure IBCC underperforms the Brok_Flw_L benchmark.

MACHINE LEARNING IN ACTION

An unhelpful aspect of machine learning systems is their reputation for being *black boxes* that users cannot understand. Whether or not one subscribes to this point of view, it is important to have easily interpreted diagnostic tools available that allow inspection of the model's internal components, especially as these evolve through time. In what follows, we provide two such tools.

Broker-Level Diagnostics

The first is an animated visualization that displays the evolution of a broker's recommendation distributions²⁴ conditional on each truth $t = 0, 1, 2$. These distributions are precisely what the system has learned about that broker's recommendation behavior up to each evaluation date. A snapshot of the animation for one broker (the one with broker code IBES_207) is given in Exhibit 15; the full animated version, which depicts the evolution of these distributions for four different brokers (IBES_199, IBES_207, IBES_410 and IBES_1296), is available online.²⁵

²⁴ Each conditional recommendation distribution is actually four-dimensional, not three-dimensional. In each case, we have marginalized over the label corresponding to Missing to obtain a three-dimensional distribution. It is these that we have plotted as triangular heatmaps.

²⁵ https://faculty.fuqua.duke.edu/~charvey/JFDS_2018/IBCC_Animation.mpeg.

EXHIBIT 9

IBCC Results with the Stock Universe Split into Two by Market Capitalization

Side	Size	Simulation Name	Return Mean	Vol	Alpha	Alpha t-Stat	Beta	Turnover
Long-Only	Large Half	Brok_Flw_L	5.95	23.93	5.96	3.74	1.01	5.79
		IBCC_Rol_L	6.96	24.57	7.22	3.01	1.01	5.75
		IBCC_Exp_L	7.70	24.50	7.98	3.44	1.01	5.71
		Both_Rol_L	7.77	24.75	7.98	3.14	1.01	6.40
		Both_Exp_L	8.11	24.51	8.38	3.43	1.00	6.32
	Small Half	Brok_Flw_L	4.70	25.60	4.76	1.63	1.03	6.00
		IBCC_Rol_L	4.15	25.98	4.54	1.73	1.06	6.08
		IBCC_Exp_L	2.42	26.08	2.76	1.06	1.07	6.03
		Both_Rol_L	5.89	25.61	6.30	2.13	1.03	6.33
		Both_Exp_L	4.66	25.78	4.98	1.72	1.04	6.27
Short-Only	Large Half	Brok_Flw_S	-2.43	24.81	-2.33	-1.08	-1.03	6.55
		IBCC_Rol_S	-5.69	25.70	-5.30	-2.93	-1.08	6.37
		IBCC_Exp_S	-5.22	25.54	-5.03	-2.46	-1.07	6.41
		Both_Rol_S	-2.53	27.90	-2.13	-0.61	-1.11	7.60
		Both_Exp_S	4.35	28.30	4.59	1.22	-1.12	7.58
	Small Half	Brok_Flw_S	1.98	27.23	2.54	0.77	-1.09	6.51
		IBCC_Rol_S	0.08	26.24	0.23	0.07	-1.05	6.54
		IBCC_Exp_S	-2.68	26.14	-2.13	-0.68	-1.05	6.51
		Both_Rol_S	8.68	27.32	9.49	2.22	-1.02	7.25
		Both_Exp_S	3.23	27.41	4.63	1.02	-1.01	7.16

Note: The alpha values are annualized.

EXHIBIT 10

Distribution of Market Capitalization after Bisecting the Universe

Universe	Mktcap Bucket	Sum Position (%)	Number of Stocks	Return (% p.a.)	Risk (% p.a.)
Large half	Mega Cap (>US\$25 billion)	41.5	75.2	4.2	6.7
	Large Cap (US\$10 billion to US\$25 billion)	33.4	70.6	3.1	6.1
	Mid Cap (US\$2 billion to US\$10 billion)	22.2	53.4	-0.3	6.9
	Small Cap (US\$250 million to US\$2 billion)	2.8	7.8	-1.4	1.9
Small half	Micro Cap (<US\$250 million)	0.0	0.1	0.0	0.1
	Mega Cap (>US\$25 billion)	0.0	0.1	0.0	0.0
	Large Cap (US\$10 billion to US\$25 billion)	5.1	9.9	1.7	0.9
	Mid Cap (US\$2 billion to US\$10 billion)	80.2	149.8	6.8	15.2
	Small Cap (US\$250 million to US\$2 billion)	14.5	30.7	-3.8	8.5
	Micro Cap (<US\$250 million)	0.1	0.2	-0.2	0.3

Notes: Here the positions for Brok_Flw_L are summarized. The numbers shown in this exhibit are time series averages 2007–2012.

The vertices of the triangles represent the three different recommendations hold, buy (here labeled Go_L), and sell (here, Go_S). Each point within a triangle corresponds to a three-vector of probabilities over

these recommendations, with the color of each point depicting its posterior probability. If the three heatmaps in Exhibit 15 were identical, then knowledge of that broker's recommendation would impart no information

EXHIBIT 11

Varying c to Change the Conviction Level Needed to Initiate a Trade for the Long-Only Models for the Period 2007–2012

Model	c	Mean	Vol	Alpha	Alpha t -Stat	Beta	Beta t -Stat	Turnover
Brok_Flw_L	—	5.43	24.18	5.47	2.73	1.01	26.97	5.75
IBCC_Exp_L	1.0	5.30	24.89	5.39	2.27	1.03	23.63	5.68
	1.1	5.13	25.38	5.43	2.06	1.04	20.67	5.81
	1.2	4.77	26.31	5.42	1.82	1.06	17.78	5.96
	1.3	6.46	26.23	6.49	2.04	1.05	16.58	6.06
	1.4	6.37	26.86	6.15	1.78	1.06	14.70	6.17
IBCC_Rol_L	1.5	5.89	27.34	5.97	1.59	1.06	13.13	6.33
	1.0	4.77	24.66	4.91	2.09	1.02	23.36	5.74
	1.1	4.82	25.14	5.12	1.98	1.03	20.62	5.83
	1.2	4.66	25.60	5.33	1.90	1.04	18.41	5.97
	1.3	5.14	26.04	5.20	1.74	1.05	16.95	5.95
Both_Exp_L	1.4	5.21	26.48	5.18	1.63	1.06	16.03	6.18
	1.5	5.55	27.01	5.56	1.61	1.07	14.56	6.38
	1.0	7.13	24.71	7.28	2.84	1.01	20.47	6.07
	1.1	7.10	25.15	7.59	2.67	1.01	18.17	6.23
	1.2	6.85	26.36	7.22	2.18	1.04	15.14	6.33
Both_Rol_L	1.3	8.73	26.67	8.90	2.50	1.04	14.26	6.49
	1.4	8.33	27.63	8.40	2.14	1.06	12.99	6.56
	1.5	8.52	28.15	8.85	2.10	1.07	11.85	6.73
	1.0	6.99	24.51	7.06	2.75	1.00	20.18	6.13
	1.1	6.98	24.86	7.40	2.63	1.00	17.93	6.28
	1.2	7.12	25.44	7.48	2.37	1.01	15.31	6.38
	1.3	7.36	26.14	7.54	2.18	1.03	13.81	6.37
	1.4	6.34	27.09	6.56	1.74	1.05	12.89	6.58
	1.5	6.97	27.76	7.12	1.76	1.06	11.91	6.78

Notes: Results are based on varying c while holding $k = 1$ in the decision rule. Recommendations were aggregated within the usual 30-day window when combining brokers. The alpha values are annualized.

about the truth outcome. In this exhibit, the three heat-maps are not identical, but the differences are subtle. This broker also displays the typical broker characteristic of having a low probability of issuing sell (Go_S) recommendations whatever the observed truth outcome.

Stock-Level Diagnostics

The focus of the previous section was visualizing what the model learns about a particular broker from the ensemble of their recommendations across a multiplicity of stocks. Here we fix our attention on a particular stock and visualize information from the multiplicity of brokers that make recommendations on that stock.

Exhibit 16 shows our visual diagnostic for the stock with identifier AST14822 (an internal code that

is unimportant). The top-panel shows the time series of recommendations for the five most prolific brokers that comment on that stock; the green and red symbols represent buy and sell, respectively, and the black symbol represents hold (labeled here as *filtered*, equivalently). The second panel lists the same information but is more cluttered because it now includes the recommendations of all brokers commenting on that stock. The third panel shows the actions that result from the predicted truths, obtained using our rolling out-of-sample process with the $c = k = 1$ case of the decision rule discussed previously. No predictions are made during the initial three-year in-sample period, so the panel is blank at the start. The fourth panel shows the positions obtained from these actions for the Both portfolios in the expanding-window long-only and long-short cases, together with

EXHIBIT 12

Varying k to Change the Conviction Level Needed to Initiate a Trade for the Long-Only Models for the Period 2007–2012

Model	k	Mean	Vol	Alpha	Alpha t -Stat	Beta	Beta t -Stat	Turnover
Brok_Flw_L	—	5.43	24.18	5.47	2.73	1.01	26.97	5.75
IBCC_Exp_L	1.0	5.30	24.89	5.39	2.27	1.03	23.63	5.68
	1.1	5.11	24.96	5.22	2.11	1.03	22.25	5.73
	1.2	5.22	25.20	5.40	2.07	1.03	20.90	5.81
	1.3	5.45	25.78	5.44	1.91	1.04	18.22	5.92
	1.4	5.30	26.25	5.22	1.74	1.06	17.36	6.16
	1.5	5.50	26.51	5.47	1.72	1.06	16.30	6.10
IBCC_Rol_L	1.0	4.77	24.66	4.91	2.09	1.02	23.36	5.74
	1.1	4.88	24.70	5.05	2.07	1.02	22.00	5.83
	1.2	4.83	24.84	5.06	1.97	1.02	20.19	5.85
	1.3	5.24	25.21	5.28	1.94	1.03	18.80	5.97
	1.4	5.23	25.57	5.21	1.83	1.04	17.73	6.12
	1.5	5.37	25.73	5.46	1.83	1.04	16.95	6.05
Both_Exp_L	1.0	7.13	24.71	7.28	2.84	1.01	20.47	6.07
	1.1	7.02	24.78	7.19	2.70	1.01	19.38	6.09
	1.2	7.24	24.89	7.48	2.69	1.01	18.64	6.17
	1.3	7.44	25.65	7.60	2.48	1.03	16.61	6.34
	1.4	7.21	26.30	7.27	2.21	1.04	15.56	6.42
	1.5	7.62	26.67	7.72	2.23	1.05	14.85	6.47
Both_Rol_L	1.0	6.99	24.51	7.06	2.75	1.00	20.18	6.13
	1.1	7.22	24.60	7.32	2.78	1.00	19.44	6.20
	1.2	7.07	24.69	7.29	2.62	1.00	18.06	6.25
	1.3	7.21	24.95	7.40	2.50	1.00	16.94	6.34
	1.4	6.35	25.36	6.48	2.06	1.01	15.54	6.37
	1.5	6.86	25.61	7.03	2.10	1.01	14.65	6.45

Notes: Results are based on varying k while holding $c = 1$ in the decision rule. Recommendations were aggregated within the usual 30-day window when combining brokers. The alpha values are annualized.

various Brok_Flw_* benchmarks. The final panel shows the truth (target) outcomes for each recommendation.

CONCLUSIONS AND FUTURE RESEARCH DIRECTIONS

We have demonstrated a computationally efficient practical approach for combining analysts' forecasts using a probabilistic machine learning model called IBCC combining it with a state-of-the-art approximate inference technique called VB. Throughout our results, the best outcomes were obtained when there was agreement between the broker recommendations and the machine learning–based forecasts obtained using IBCC. These findings echo important current research in the area of human–computer interaction, where decision

making based on inputs from artificial intelligence and other sources is used to assist human decision making. It also suggests some intriguing research directions for enhancing the investment processes and performance of both quantitative and discretionary fund managers.

An important advantage of the IBCC model is its scalability compared to other multivariate dependence techniques (e.g., copula models). Our application integrated recommendations from 347 brokers; however, IBCC has been successfully used in applications involving many thousands of individual classifiers, so there is ample scope for extension. For example, we could look at individual analysts, or more refined groups of analysts, rather than brokers.²⁶ In addition, it may

²⁶We are unable to report on this because of current restrictions.

EXHIBIT 13

Varying the Truth Boundary Parameter over the Period 2007–2012 for Long-Only and Short-Only Portfolios

Model	Truth (%)	Mean	Vol	Alpha	Alpha t-Stat	Beta	Beta t-Stat	Turnover
Brok_Flw_L	—	5.43	24.18	5.47	2.73	1.01	26.97	5.75
Both_Exp_L	1	6.07	24.70	6.04	2.55	1.02	21.53	6.00
	2	6.09	24.72	6.24	2.59	1.02	21.34	6.03
	3	6.39	24.65	6.56	2.68	1.01	21.24	6.00
	4	6.57	24.63	6.70	2.68	1.01	20.93	6.03
	5	7.13	24.71	7.28	2.84	1.01	20.47	6.07
	6	7.61	24.86	7.81	2.90	1.01	19.82	6.25
	8	5.94	26.31	6.24	1.86	1.04	15.25	6.56
	10	6.66	26.95	6.74	1.84	1.05	14.93	6.77
Both_Rol_L	1	6.29	24.39	6.30	2.66	1.00	21.82	6.00
	2	6.10	24.39	6.26	2.62	1.00	21.41	6.03
	3	6.10	24.42	6.26	2.59	1.00	21.34	6.00
	4	6.32	24.43	6.46	2.61	1.00	20.69	6.05
	5	6.99	24.51	7.06	2.75	1.00	20.18	6.13
	6	6.97	24.60	7.29	2.75	1.00	19.98	6.29
	8	7.17	26.19	7.95	2.33	1.03	14.78	6.67
	10	7.74	27.00	7.59	2.13	1.06	15.09	6.87
Brok_Flw_S	—	−0.51	24.96	−0.13	−0.05	−1.03	−29.42	6.38
Both_Exp_S	1	2.10	25.76	2.31	0.83	−1.05	−26.19	6.91
	2	3.36	25.56	3.55	1.27	−1.04	−26.51	6.94
	3	3.13	25.45	3.38	1.23	−1.04	−26.27	6.97
	4	2.17	25.53	2.35	0.81	−1.03	−25.31	7.00
	5	2.11	25.71	2.46	0.79	−1.03	−25.39	7.04
	6	4.09	26.51	4.02	1.05	−1.02	−19.30	7.20
	8	0.67	28.84	1.06	0.21	−1.04	−15.77	7.50
	10	−1.15	31.44	−1.28	−0.22	−1.09	−14.66	7.76
Both_Rol_S	1	2.73	25.68	2.92	1.09	−1.05	−25.52	6.79
	2	2.75	25.43	2.93	1.08	−1.04	−24.57	6.78
	3	3.28	25.51	3.49	1.27	−1.04	−24.65	6.84
	4	1.80	25.69	1.98	0.69	−1.04	−24.71	6.99
	5	2.99	25.98	3.45	1.06	−1.03	−21.81	7.12
	6	2.71	26.76	3.07	0.88	−1.05	−23.33	7.23
	8	0.36	27.63	−0.24	−0.06	−1.05	−20.90	7.50
	10	−1.18	29.56	−1.81	−0.36	−1.08	−20.01	7.62

Notes: Unreported results show that IBCC_* consistently underperforms Brok_Flw and Both_*, consistent with our previous findings. The alpha values are annualized.

be useful to combine the recommendation data examined here with categorical sentiment measures extracted using a range of different natural language interpreters on both mainstream and financial news sources. There is scope to obtain an order of magnitude more classifiers. The computational efficiency of our implementation would enable such data to be handled without issue and real time forecasting to be undertaken.

Although the VB implementation of the IBCC holds promise, it also has limitations. In the Galaxy Zoo experiment that we used to motivate the research application, several distinct issues make the application to analysts different from the application to astronomers. First, it is reasonable to assume that the astronomers are operating independently (not collaborating with each other). However, it is likely that analysts are aware of

EXHIBIT 14

Varying the Holding Period for the Long-Only Models for Period 2007–2012

Model	Holding Period	Mean	Vol	Alpha	Alpha <i>t</i> -Stat	Beta	Beta <i>t</i> -Stat	Turnover
Brok_Flw_L	10	10.04	23.99	10.01	4.55	0.99	29.08	28.20
	20	7.39	24.16	7.27	3.51	1.01	28.17	14.89
	30	6.82	24.33	6.71	3.23	1.02	26.99	10.35
	45	6.05	24.27	6.01	2.98	1.01	26.48	7.31
	60	5.43	24.18	5.47	2.73	1.01	26.97	5.75
	90	5.46	23.94	5.32	2.75	1.00	29.61	4.31
IBCC_Exp_L	10	7.09	24.01	7.08	2.92	0.98	25.66	28.73
	20	5.93	24.81	5.76	2.29	1.02	21.80	15.19
	30	5.15	24.77	5.18	2.13	1.02	21.93	10.55
	45	5.61	24.79	5.76	2.43	1.02	23.20	7.29
	60	5.30	24.89	5.39	2.27	1.03	23.63	5.68
	90	4.95	24.62	5.10	2.23	1.02	26.34	4.20
IBCC_Rol_L	10	8.02	24.01	7.84	3.15	0.98	22.83	28.76
	20	5.75	24.76	5.64	2.22	1.01	21.09	15.24
	30	6.13	24.55	6.20	2.56	1.01	22.49	10.53
	45	5.60	24.46	5.79	2.52	1.01	24.75	7.31
	60	4.77	24.66	4.91	2.09	1.02	23.36	5.74
	90	5.12	24.58	5.28	2.31	1.02	25.71	4.21
Both_Exp_L	10	14.05	23.97	14.21	5.51	0.97	25.38	30.18
	20	9.76	24.75	9.80	3.76	1.01	20.78	15.99
	30	8.42	24.79	8.24	3.18	1.01	20.87	11.03
	45	7.98	24.91	7.93	3.04	1.02	20.73	7.74
	60	7.13	24.71	7.28	2.84	1.01	20.47	6.07
	90	6.63	24.39	6.65	2.66	1.00	22.05	4.71
Both_Rol_L	10	13.38	23.87	13.36	5.00	0.96	22.10	30.28
	20	9.44	24.35	9.62	3.63	0.99	20.48	16.06
	30	9.27	24.66	9.11	3.41	1.00	21.07	11.03
	45	7.71	24.51	7.75	3.07	1.00	21.43	7.73
	60	6.99	24.51	7.06	2.75	1.00	20.18	6.13
	90	6.48	24.27	6.61	2.66	0.99	21.83	4.67

Notes: Here the trade holding period and the holding period for assessing truths are constrained to be equal. The $\pm 5\%$ threshold for converting stock returns to truths is scaled to yield a similar number of truths for each horizon, using the usual random walk property that $\sigma(X_t) \propto \sqrt{t}$, which gives threshold $= \sqrt{t/t_0} \times 5\% = \sqrt{\text{holding period}/60} \times 5\%$. As elsewhere, recommendations are aggregated with a lookback of up to 30 days when combining brokers. The alpha values are annualized.

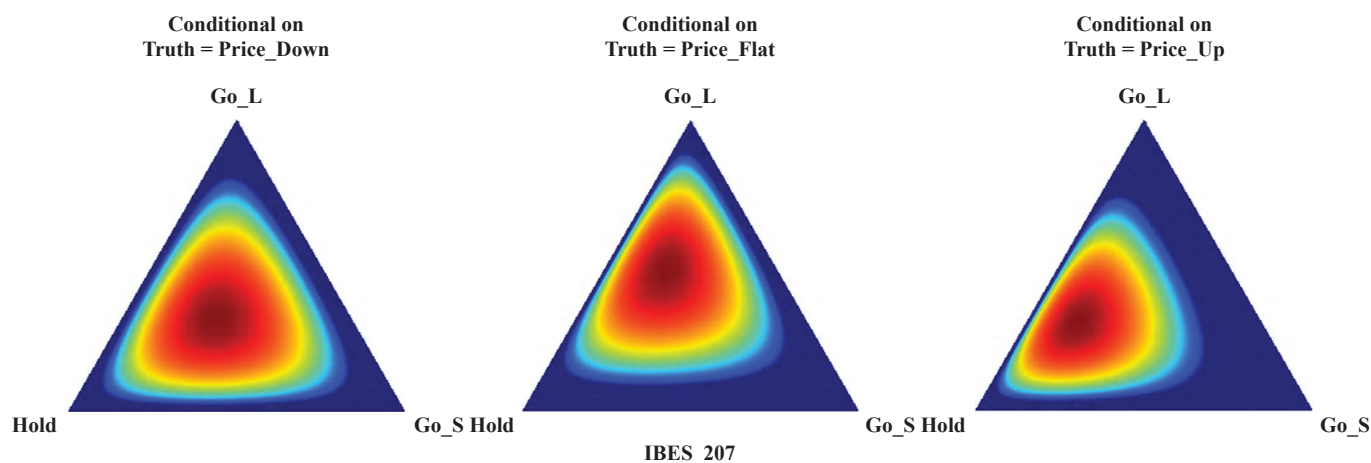
other analysts' forecasts—and this could affect their forecasts. Second, the quality of all analysts' forecasts could be affected by common events such as a recession, global sentiment, and common market factors that may affect their sector or region. Such common factors do not apply in the Galaxy Zoo experiment.

There are also areas for methodological consideration within the current implementation. For example, the IBCC model has no concept of ordering within the truth outcomes or the recommendations; they are

simply sets of categorical labels. Perhaps more importantly, IBCC has no concept of parity between the recommendations and truths. Maybe it is therefore only to be expected that our strongest results arose when we looked for reinforcement between the raw broker recommendations and the IBCC predictions. Changing the model to incorporate some parity effect would make it less general but would likely boost performance in our application. On the other hand, if sufficient data were available to learn the parity relationship with the

EXHIBIT 15

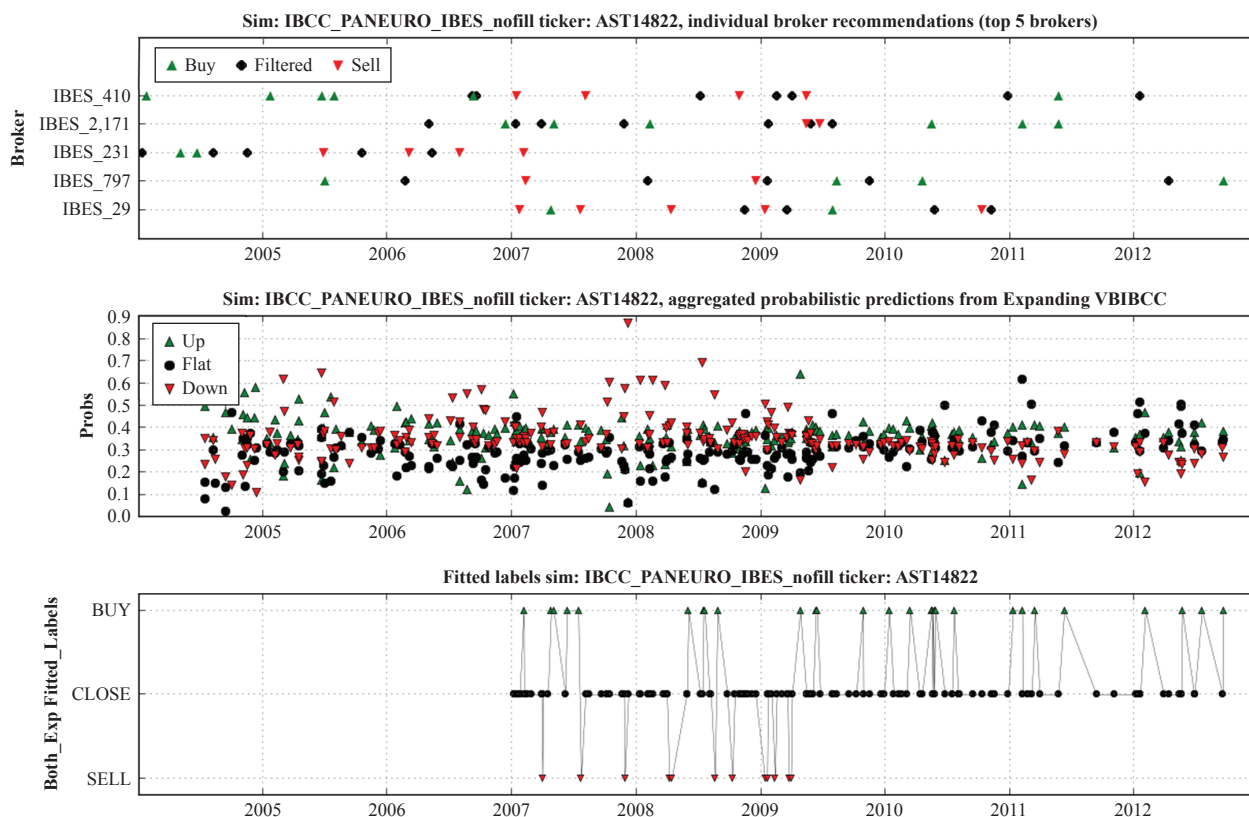
Screenshot of the Website Animation



Notes: This exhibit shows the evolving recommendation distributions for the broker with identifier IBES_207, conditional on truth=Price_Down (left), truth=Price_Flat (middle), and truth=Price_Up (right). Within each triangle, each pixel represents a three-vector of probabilities over the recommendations hold, buy (Go_L), and sell (Go_S). The color of each pixel represents the posterior probability of this corresponding three-vector; blue pixels have very low posterior probability, and dark red pixels have the highest posterior probability.

EXHIBIT 16

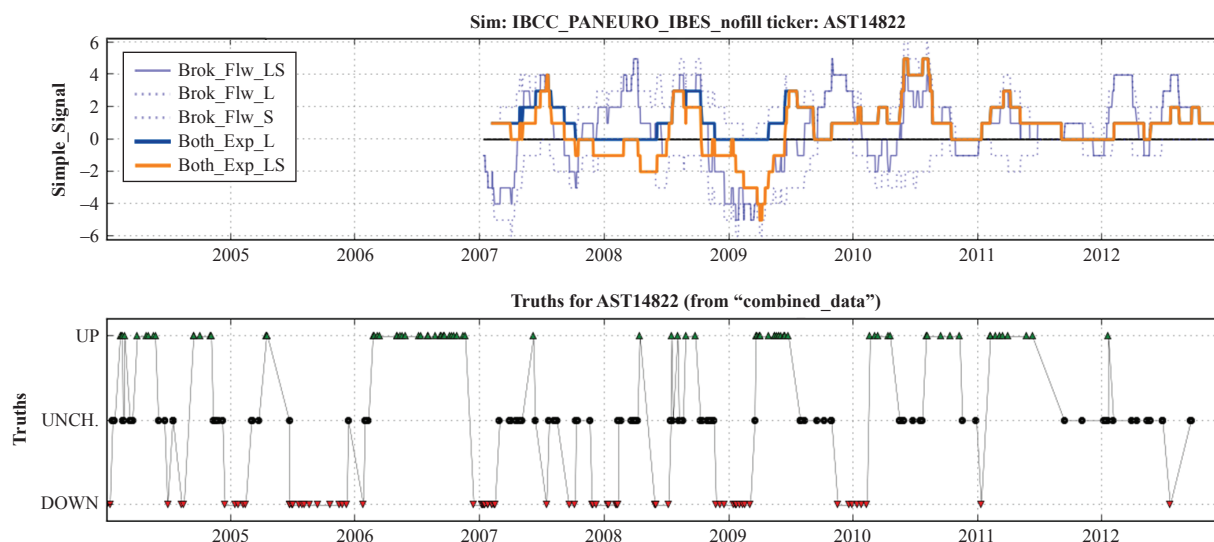
Diagnostic Panel for the Stock with Identifier AST14822



(continued)

EXHIBIT 16 (continued)

Diagnostic Panel for the Stock with Identifier AST14822



Notes: The top panel shows recommendations for the five most prolific brokers that comment on the stock; green and red represent buy and sell, respectively, and the black symbol represents hold (labeled here as filtered, equivalently). The second panel lists the same information and now includes the recommendations of all brokers commenting on that stock. The third panel shows the actions that result from the predicted truths, obtained using our rolling out-of-sample process with the $c = k = 1$ case of the decision rule. No predictions are made during the initial three-year in-sample period. The fourth panel shows the positions obtained from these actions for the Both portfolios in the expanding-window long-only and long-short cases, together with various *Brok_Flw_** benchmarks. The final panel shows the truth (target) outcomes for each recommendation.

original IBCC model, then there would be no issue. Our practical experience is that there are never sufficient data available compared to what we would like, so working with flexible but not completely general models gives the best results.

ACKNOWLEDGMENT

We thank Edwin Simpson for supporting discussions.

REFERENCES

- Bernardo, J. M., and A. F. M. Smith. *Bayesian Theory*. Hoboken: Wiley, 1994.
- Bishop, C. *Pattern Recognition and Machine Learning*. New York: Springer, 2006.
- Blei, D. M., A. Kucukelbir, and J. D. McAuliffe. 2018. "Variational Inference: A Review for Statisticians." arXiv:1601.00670v9.
- Boyd, S., and L. Vandenberghe. *Convex Optimization*. Cambridge: Cambridge University Press, 2004.
- Bradshaw, M. T. 2011. "Analysts' Forecasts: What Do We Know after Decades of Work?" SSRN, June 30, <https://ssrn.com/abstract=1880339>.
- Brown, L. 1993. "Earnings Forecasting Research: Its Implications for Capital Markets Research." *International Journal of Forecasting* 9: 295–320.
- , ed. 2000. *I/B/E/S Research Bibliography*. 6th ed. New York: I/B/E/S International Inc., 2000.
- Dempster, A. P., N. M. Laird, and D. B. Rubin. 1977. "Maximum Likelihood from Incomplete Data via the EM Algorithm." *Journal of the Royal Statistical Society, Series B (Methodological)* 39 (1): 1–38.
- Fox, C. W., and S. J. Roberts. 2011. "A Tutorial on Variational Bayesian Inference." *Artificial Intelligence Review* 38 (2): 85–95.

- Ghahramani, Z., and H. C. Kim. 2003. "Bayesian Classifier Combination." Gatsby Computational Neuroscience Unit technical report no. GCNU-T. London, UK.
- Givoly, D., and J. Lakonishok. 1984. "Properties of Analysts' Forecasts of Earnings: A Review and Analysis of the Research." *Journal of Accounting Literature* 3: 117–152.
- Kim, H. C., and Z. Ghahramani. 2012. "Bayesian Classifier Combination." Proceedings of the 15th AISTATS Conference.
- Lee, P. M. *Bayesian Statistics: An Introduction*. Chichester, UK: John Wiley, 2012.
- Levenberg, A., S. Pulman, K. Moilanen, E. Simpson, and S. Roberts. 2014. "Predicting Economic Indicators from Web Text Using Sentiment Composition." *International Journal of Computer and Communication Engineering* 3 (2): 109–115.
- Levenberg, A., E. Simpson, S. Roberts, and G. Gottlob. "Economic Prediction Using Heterogeneous Data Streams from the World Wide Web." In *Scalable Decision Making: Uncertainty, Imperfection, Deliberation (SCALE)*, Proceedings of ECML/PKDD Workshop. New York: Springer, 2013.
- Lintott, C. 2012. "I, for One, Welcome Our New Machine Collaborators." August 3, <https://blog.zooniverse.org/2012/08/03/i-for-one-welcome-our-new-machine-collaborators>.
- Parisi, G. *Statistical Field Theory*. Boston: Addison-Wesley, 1988.
- Sato, M. A. 2001. "Online Model Selection Based on the Variational Bayes." *Neural Computation* 13: 1649–1681.
- Schipper, K. 1991. "Commentary on Analysts' Forecasts." *Accounting Horizons* 5 (4): 105–121.
- Simpson, E., S. Roberts, I. Psorakis, and A. Smith. "Dynamic Bayesian Combination of Multiple Imperfect Classifiers." In *Decision Making and Imperfection. Intelligent Systems Reference Library Series*, Vol. 474. New York: Springer, 2013.
- Tanner, M. A. *Tools for Statistical Inference: Methods for the Exploration of Posterior Distributions and Likelihood Functions*. New York: Springer, 1996.
- To order reprints of this article, please contact David Rowe at d.rowe@pageantmedia.com or 646-891-2157.