

Frontiers: News Event-Driven Forecasting of Commodity Prices

Sunandan Chakraborty,^a Srikanth Jagabathula,^b Lakshminarayanan Subramanian,^c Ashwin Venkataraman^{d,*}

^aLuddy School of Informatics, Computing and Engineering, Indiana University, Indianapolis, Indiana 46202; ^bStern School of Business, New York University, New York, New York 10012; ^cCourant Institute of Mathematical Sciences, New York University, New York, New York 10012; ^dNaveen Jindal School of Management, University of Texas at Dallas, Richardson, Texas 75080

*Corresponding author

Contact: sunchak@iu.edu,  <https://orcid.org/0000-0002-3331-6082> (SC); sjagabat@stern.nyu.edu,  <https://orcid.org/0000-0002-4854-3181> (S); lakshmi@cs.nyu.edu (LS); axv190029@utdallas.edu,  <https://orcid.org/0000-0002-6182-2361> (AV)

Received: December 21, 2022

Revised: June 1, 2023; October 31, 2023; January 5, 2024

Accepted: January 13, 2024

Published Online in *Articles in Advance*: March 20, 2024

<https://doi.org/10.1287/msom.2022.0641>

Copyright: © 2024 INFORMS

Abstract. *Problem definition:* Commodity prices have exhibited significant volatility in recent times, which poses an exogenous risk factor for commodity-processing and commodity-trading firms. Accurate commodity price forecasts can help firms leverage data-driven procurement policies that incorporate the underlying price volatility for financial and operational hedging decisions. However, historical prices alone are insufficient to obtain reasonable forecasts because of the extreme volatility. *Methodology/results:* Building on the hypothesis that commodity prices are driven by real-world events, we propose a method that automatically extracts events from news articles and combines them with price data using a neural network-based predictive model to forecast prices. In addition to achieving a high prediction accuracy that outperforms several benchmarks (by up to 13%), our proposed model is also *interpretable*, which allows us to identify meaningful events driving the price fluctuations. We found that the events frequently associated with major fluctuations in the price include “natural,” “hike,” “policy,” and “elections,” all of which are known drivers of price change. We used a corpus containing about 1.6 million news articles of a major Indian newspaper spanning 15 years and daily prices of four crops (onion, potato, rice, and wheat) in India to perform this study. Our proposed approach is flexible and can be used to predict other time series data, such as disease incidence levels or macroeconomic indicators, that are also influenced by real-world events. *Managerial implications:* Firms can leverage price forecasts from our system to design inventory and procurement policies in the face of uncertain commodity prices. Commodity merchants can also use the forecasts to design optimal storage policies for physical trading of commodities when prices are volatile. Our findings can also significantly impact policymakers, who can leverage the information of impending price changes and associated events to mitigate the negative effects of price shocks.

History: This paper has been accepted in the *Manufacturing & Service Operations Management* Frontiers in Operations Initiative.

Supplemental Material: The online appendix is available at <https://doi.org/10.1287/msom.2022.0641>.

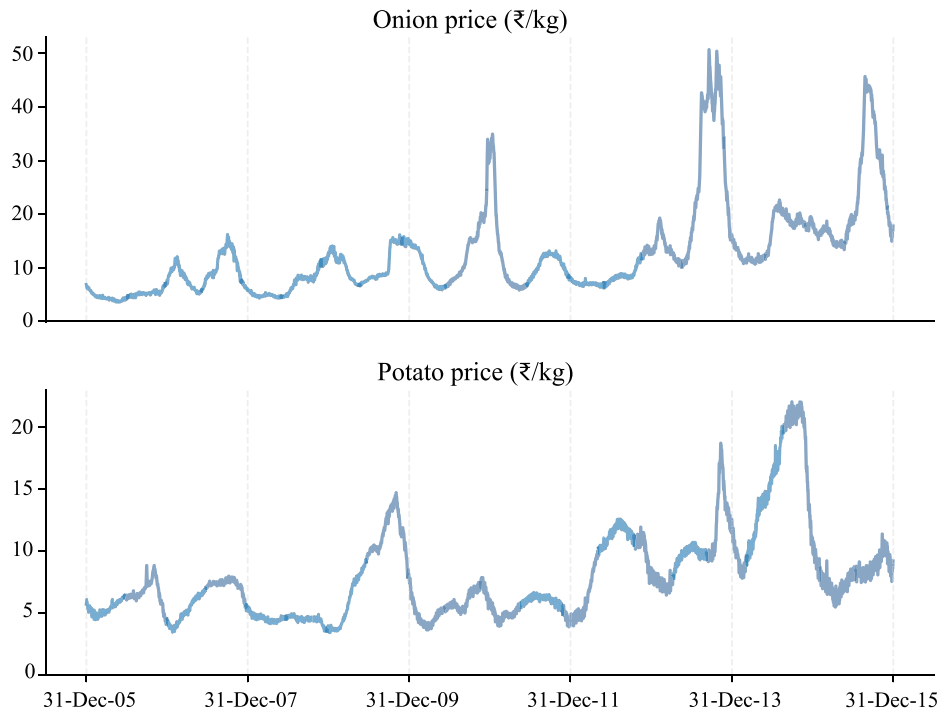
Keywords: forecasting • price volatility • food supply chain • textual data • procurement

1. Introduction

With the rise of supply chain analytics, classic operational decisions within supply chains, such as sourcing and procurement, are increasingly becoming data driven (Mandl 2019). These decisions critically rely on accurate forecasts of not just the downstream demand but also upstream costs. For example, a significant fraction of the total costs of various firms is composed of commodity costs: 27% in the mechanical and plant engineering industries, 47% in the automotive supply industry, 56% in the packaging industry, and 66% in the agri-food industry (Mandl 2019). Therefore, volatility in commodity prices significantly affects the procurement costs of these firms and poses an exogenous

risk factor in practice. Accurate commodity price forecasts allow firms to use data-driven procurement policies that incorporate the underlying price volatility to hedge against price risk (Mandl and Minner 2023) and effectively control inventory (Berling and Martínez-de Albéniz 2011, Goel and Gutierrez 2011, Xiao et al. 2015), which can be viewed as a form of operational hedging. Extant literature in operations has mostly focused on downstream demand forecasting (see, e.g., Kurawarwala and Matsuo 1996, Carbonneau et al. 2008, Boone et al. 2018), but despite its significance, there is a dearth of OM literature on forecasting purchasing costs or other supply-side factors (Haksöz and Seshadri 2007, Syntetos et al. 2016).

Figure 1. (Color online) Daily National Average Onion (Upper Panel) and Potato (Lower Panel) Prices in India from 2006 to 2015



In this paper, we consider the problem of forecasting agricultural commodity prices in India, one of the largest economies in the world. In addition to its relevance to the procurement decisions mentioned earlier, this is a hugely important problem because India is still primarily an agrarian economy, where agriculture-related sectors (such as forestry and fisheries) are responsible for over 50% of its labor force and account for 20% of the country's GDP (International Trade Administration 2022). Indeed, prior studies have shown that rising food prices can reduce household welfare (Mahajan et al. 2015, Weber 2015), affect livelihoods of farmers and rural populations (De Janvry and Sadoulet 2009, Pons 2011), increase food insecurity and human development vulnerabilities of women and children (Dev 2011), indirectly increase crime rates (New York Times 2013, India.com 2019), increase illegal hoarding (Sharma et al. 2011, Levi et al. 2022), and even influence election outcomes (McLain 2013, Virk 2015, Agarwal 2019, CNBC-TV18 2019, Deuskar 2020). Mitigating these negative effects requires effective policy design that relies on not just anticipating the price fluctuations but also understanding the drivers behind sudden shocks (Saha et al. 2020, Saxena et al. 2020).

Predicting the fluctuations, however, often requires more than just the past price data. These data are inherently volatile without discernible patterns. For example, the national average onion price in India increased more than 400% within a span of seven days in December

2012; such abrupt changes are surprisingly frequent and generally do not follow any seasonal trends as seen in Figure 1. Such extreme volatility is not restricted to the Indian context. For instance, figure 2.1 in Mandl (2019) shows that a wide array of commodities, including metals (silver, lead, etc.), agricultural products (corn, wheat, coffee, etc.), and energy (gas, oil, etc.), has been significantly more volatile than exchange rates and corporate stocks over the past two decades.

More importantly, historical price data lack information on several other exogenous factors, such as the economic climate, temperature, inventories, interest, and/or exchange rates, which impact prices (see, e.g., Pindyck 2004). In the context of India, numerous prior studies (see, e.g., Headey and Fan 2008, Raka and Ramesh 2017, Birthal et al. 2019) have identified a number of causal factors for food price rise, including nonseasonal rainfall in the production belt; lack of rainfall during cultivation; and natural disasters, such as floods, that affect both production and supply. Prices are also routinely affected by local transport strikes, fluctuations in fuel prices, and festivals. To account for these drivers, existing methods (see, e.g., Madaan et al. 2019) have mostly relied on *structured* data sources and manual selection of a predetermined set of driving factors for predicting food price fluctuations. Such structured approaches improve prediction accuracy, but manual preselection can easily miss out on several key factors. Prices may rise now because of

floods but later because of a sudden tax hike, both of which may not be captured together in a single structured data source. To meaningfully improve our prediction accuracy, we need an automatic and scalable way to identify and account for the many factors, which tend to evolve, merge, and disappear over time.

Our approach to price prediction relies on moving away from structured data sources to automatically extracting the relevant factors from a rich source of *unstructured data*—online news articles. News streams provide us with high-frequency data on events or factors that tend to affect food prices. For example, Figure 2(a) shows the price fluctuations on the daily average onion price in India during a one-year window from November 2008 to October 2009 associated with a small sample set of news article headlines reporting the fluctuation itself or the cause behind the fluctuation. Examples of events responsible for onion price variations captured in Figure 2 include crop infection at the end of 2008 (TOI 2008) further worsened by transport strikes and an oil crisis at the beginning of 2009 (TOI 2009a) followed by a price decrease triggered by influx of supply (TOI 2009b) and the latter half of 2009 exhibiting a sudden price hike triggered by excessive rains and floods (TOI 2009c, d). Figure 2(b) shows a similar plot for potato prices.

1.1. Contributions and Positioning in the OM Literature

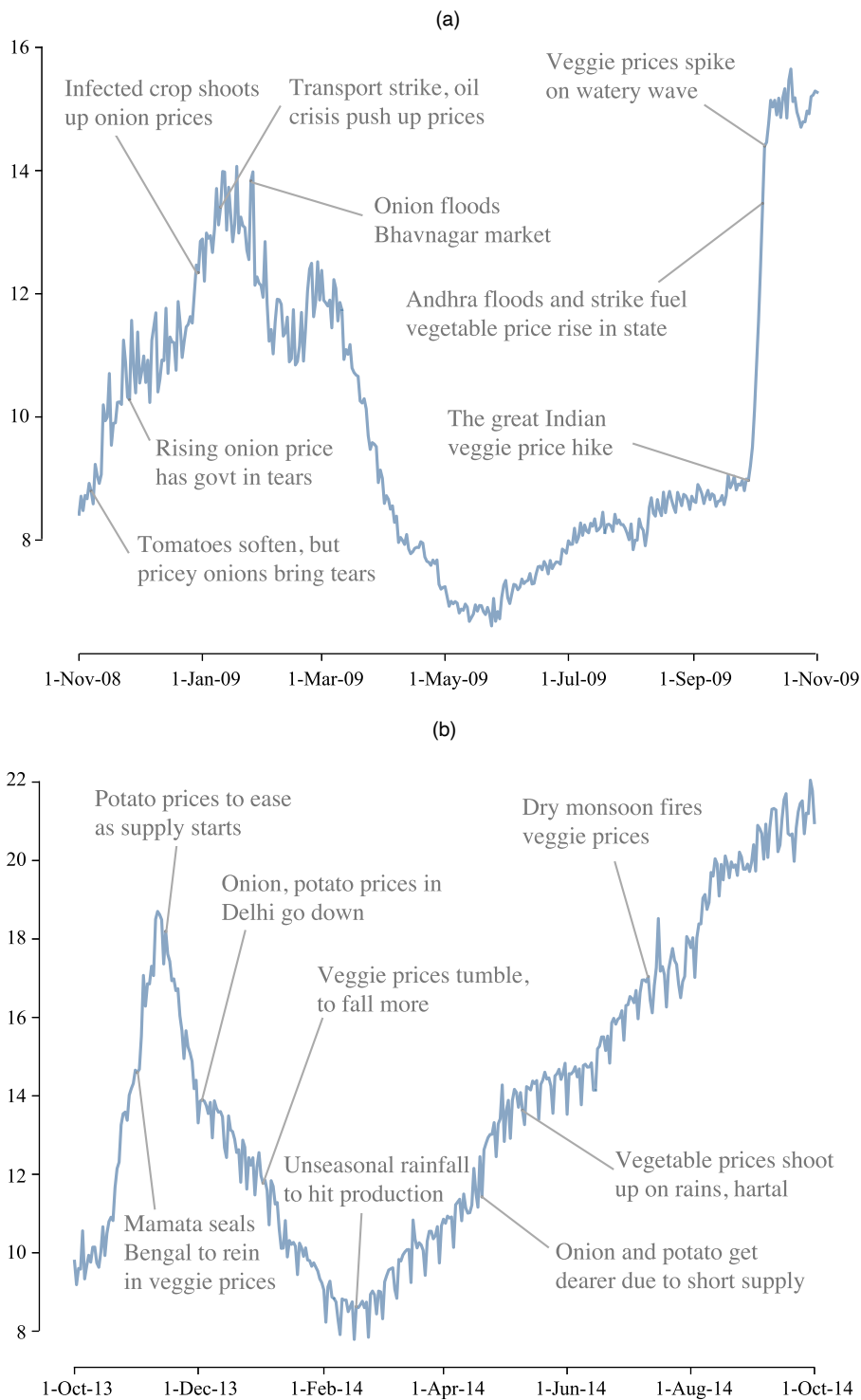
We propose an event-based methodology, which relies on historical prices and text of news articles, to forecast future commodity prices. Our method involves two key steps: (a) automatically extracting meaningful “events” from a corpus of millions of news articles and (b) relating past changes in the event intensities to past fluctuations in commodity prices to effectively forecast future prices using observed events. For the first step, we propose a scalable method that transforms each news article into what we call its “event representation,” which is a low-dimensional vector with each element indicating the intensity of the corresponding event in the article. For the second step, we design what we call a *recurrent event network* (REN), a unified deep learning architecture that jointly models the event representations and past values of commodity prices to forecast future prices. We note that a preliminary version of this work appeared in Chakraborty et al. (2016).

An important contribution of our work is a customized method to automatically extract the key events reported in a news article. It is critical for us to extract *meaningful* and *relevant* events from the text of a news article. Because a news article mentions many facts beyond just the events, such as location, entities involved, etc., general-purpose representation schemes might identify “topics” or themes that are not directly related to the events being reported. As a result, these

representations can and do indeed worsen the prediction performance (see Section 4.2). Instead, our approach relies on exploiting the specific structure of the news articles by (a) restricting the vocabulary to only consider “event triggers,” which are words/phrases that indicate the type and occurrence of an event, and (b) inferring the events reported by leveraging the unique discourse structure in news articles. Using two crowdsourcing studies, we demonstrate that our method scales to millions of news articles and automatically identifies meaningful events. See Section 3.2 for the precise description of the embedding scheme and details of the crowdsourcing studies.

The second key contribution of our work is the REN model, which uses the event embeddings as exogenous features to make a more informed prediction of the outcome variable. Unlike traditional time-series forecasting models like ARIMA (Brockwell and Davis 2002), RENs are designed to capture nonlinear dependencies between the input and output variables. They employ a recurrent neural network (RNN) architecture to model the sequential dependence of past prices and events on future prices (that is, today’s price depends on a history of past prices and events) using hidden states that propagate this dependence. However, the standard RNN architectures (sometimes called a “vanilla RNN”) cannot effectively capture long-term dependencies, such as the impact of price and events a week ago on today’s price. Because many real-world events can exhibit a delayed and long-term impact on commodity prices, we have designed RENs to employ an RNN with long short-term memory (LSTM) hidden states, which have been shown to better capture long-term dependencies. These design choices make an REN a unique data integration model, where it can seamlessly integrate structured (e.g., time-series data) and unstructured (e.g., news text) data; see Section 3.3 for the precise description of RENs.

We leverage the REN framework to forecast prices of four essential staple crops—onion, potato, rice, and wheat—in the Indian context using real-world events extracted from online news streams. Rice and wheat are consumed on a daily basis by a large fraction of the Indian population (Mittal et al. 2018). Onion and potato are among the most used vegetables in a large number of Indian dishes, and the high volatility of their prices has had a widespread and profound impact across the society (Paul et al. 2015, Gro Intelligence 2016). Our evaluation spanning a 15-year period comprising around 1.6 million news articles shows that RENs outperform several benchmark methods by up to 7% for onion and potato, 13% for rice, and 5% for wheat in one-day-ahead forecast accuracy. The improved accuracy can benefit farmers and traders in India, who typically rely on price forecasts for their buying and selling

Figure 2. (Color online) News Headlines Correlated with Fluctuations in Crop Prices

Notes. We see evidence of factors (e.g., rainfall, floods, strikes, etc.) influencing crop price fluctuations reported in news articles, which motivates our approach of leveraging news articles to improve crop price forecasting. (a) Daily onion prices (rupees per kilogram) between November 2008 and November 2009. (b) Daily potato prices (rupees per kilogram) between October 2013 and October 2014.

decisions (Digital Green 2020). The results also showcase the value of extracting the “right” set of features from text data because our custom event embeddings

that leverage the special structure of news articles outperform general-purpose embedding techniques, such as word2vec (Mikolov et al. 2013), doc2vec (Le and

Mikolov 2014), and latent dirichlet allocation (LDA) topics (Blei et al. 2003), by up to 20%. See Section 4 for more details on the numerical evaluation.

Although firms and policymakers can incorporate the knowledge of impending price fluctuations in their decision making, they are often also interested in understanding the key drivers behind the fluctuations. Because they are trained on news events, RENs naturally provide interpretability to the price forecasts. Specifically, an REN is capable of determining the relative importance of each event in influencing the price of a crop during a specific time interval. Using these event importance scores, RENs can provide a fine-grained estimate of the most important events behind specific episodes of large price fluctuations; see Section 5 for the details. We find that our method automatically extracts event classes, such as “natural” (floods, cold waves, poor rainfall, etc.), “hike” (hikes in fuel prices, fares, highway tolls, etc.), “policy” (government actions curbing/promoting imports and/or exports, etc.), and “elections,” all of which are known drivers of price changes. We also find that the key events driving price fluctuations for a particular crop change over time, and different events may drive the fluctuations in different crops, even during the same time. These findings highlight the limitations of a structured approach relying on a predefined set of events, which might miss out on relevant factors or incorrectly attribute certain events to the changing prices. Instead, our automated method relying on unstructured data casts a wide net and is less likely to miss out on key factors.

Although this paper focuses on forecasting prices of agricultural commodities, our methodology can be readily used to forecast prices of other commodities, such as metals (e.g., aluminum, copper, etc.) and energy (e.g., electricity, natural gas, etc.). As discussed earlier, such price forecasts can then be leveraged to design inventory and procurement policies that are aware of the underlying price volatility. In addition, they can be used to design optimal storage policies for physical trading of commodities under volatile prices (Secomandi 2015, Mandl et al. 2022). In fact, the REN framework proposed in this paper is very flexible and can be applied to a variety of other domains as well. As a concrete illustration, we leverage our method to forecast the number of reported cases of three infectious diseases prevalent in India. This was a critical problem during the recent COVID-19 pandemic because accurate prediction of future cases is extremely valuable for decision making at multiple levels: from healthcare providers (staffing, hospital bed allocation, etc.) to governments (awareness campaigns, resource allocation, etc.) to pharmaceutical companies (production schedule of vaccines and drugs). Our results show that RENs outperform the

benchmark models, achieving mean absolute percentage error (MAPE) values of 33.9%, 3.2%, and 13.6% for forecasting monthly reported cases of malaria, dengue, and influenza, respectively. Refer to Online Appendix F for the details.

Our work adds to the literature on developing data-driven methodologies to improve forecasting in the operations literature. However, unlike existing works that assume access to structured data sources (Ban and Rudin 2019, Cortazar et al. 2019, Zhu et al. 2021, Mandl and Minner 2023), our work makes use of unstructured text data. To the best of our knowledge, there is very limited work on utilizing text data in the OM literature—see Cui et al. (2018) and Wu (2023) for notable exceptions—and we believe our paper helps to demonstrate the rich potential of unstructured data to improve operational decisions.

2. Related Literature

We discuss the most relevant streams of literature from a methodological standpoint.

2.1. Text Mining on News Articles and Social Media

Advances in text mining techniques over the past two decades have made it possible to extract knowledge from vast unstructured text corpora in an automated manner. In particular, there has been a lot of interest in extracting structure from news articles and social media platforms using text mining techniques. Shahaf and Guestrin (2010) developed a principled approach to connect news articles related to a common event or topic to enable better understanding of news stories. Bagozzi and Schrodtt (2012) apply LDA on a vast corpus of political news reports and study the overlap between the recovered topics and those represented in existing event ontologies. The authors show that the topic modeling approach is able to uncover events that were not captured by the ontologies, suggesting the need for more automated event encoding mechanisms. Trend analysis model (Kawamae 2011) and temporal-LDA (Wang et al. 2012) model the temporal aspect of topics in social media streams, like Twitter. Vaca et al. (2014) used a collective matrix factorization method to track emerging, fading, and evolving topics in news streams. In most of these works, events and/or topics have just been used as a tool for knowledge acquisition or information extraction, whereas our goal is to use the extracted events to predict fluctuations in commodity prices.

2.2. Event Extraction in Natural Language Processing

The natural language processing (NLP) literature has focused extensively on defining and extracting events

from text data, such as news articles and Twitter streams; see Hogenboom et al. (2011) and Atefeh and Khreich (2015) for detailed surveys. It will be impossible to summarize all of this work; we only discuss a few papers here that are closest to our work. Chambers and Jurafsky (2011) use LDA and hierarchical clustering to group event templates defined as consisting of a role (typically the participants) and an event pattern (the set of words/phrases that describe the event). However, they specifically focused on events relating to terrorism. Chambers (2013) proposes a similar generative model as ours but focuses on named entities and coreference resolution of the different mentions of the same entity. Cheung et al. (2013) propose a probabilistic model for identifying semantic frames, which are a group of related events. This is similar to the notion of an event class in our work, but Cheung et al. (2013) also model the entities involved, whereas we only focus on the events. In particular, these works focus on identifying events at a sentence level and therefore, exploit very local structure, such as extracting the verb clauses in each sentence. Our work, on the other hand, is focused on identifying the most salient events that are reported anywhere in a news article. More recently, Caselli and Vossen (2017) introduce The Event Story-Line corpus as a benchmark for research aimed at identifying causal or temporal relations between events. Balashankar et al. (2019) aims to uncover *causal* relationships between events spread across time in a corpus of news articles using the notion of Granger causality. We do not aim to identify temporal or causal links between events in different time periods; rather, the focus is on minimizing the forecast error, and any dependence between events is learned implicitly by the neural network in our REN framework.

2.3. Predicting Real-World Indicators Using Text Data

There is significant research on leveraging text data to predict real-world indicators, and our work falls into this literature. Much of this work has focused on forecasting stock prices and financial indicators (see Nasirtoussi et al. 2014 and Xing et al. 2018 for detailed surveys), but other applications include forecasting economic trends and macroeconomic variables, such as GDP (Larsen and Thorsrud 2019, Bybee et al. 2021) or disease outbreaks (Chakraborty and Subramanian 2016). To keep the discussion focused, we only review literature that leverages news text.

Existing work has proposed different ways to extract information from news text. The earliest work in this stream was by Gidofalvi (2001), who models individual words in financial news articles using a naive Bayes text classifier to predict the movement of stock prices. Schumaker and Chen (2009) compared the performance of different representation schemes, such as bag of words,

noun phrases, and named entities, for improving stock price forecasting. Hagenau et al. (2013) demonstrate the value of capturing word combinations instead of single words and market feedback to select relevant features to better forecast stock prices. Other work has explored the use of NLP techniques, such as sentiment analysis (Tetlock 2007, Zhang and Skiena 2010, Li et al. 2014) and semantic frames (Xie et al. 2013), to extract relevant features from news text. More recently, machine learning (ML) approaches, such as matrix factorization (Ming et al. 2014) and neural network embeddings (Hu et al. 2018, Liu et al. 2018), have been utilized to further improve the quality of the extracted text representations. The aforementioned works apply or extend representation schemes developed in the NLP/ML literature to extract information from *generic* documents, not accounting for the particular structure of a news article (title, lead paragraph, etc.). In contrast, our proposed approach exploits the typical discourse structure of news articles to define a customized embedding that outperforms general-purpose embeddings.

There is also some work that explicitly relies on extracting some notion of events from news articles. For instance, Ding et al. (2014) extract event representations as a tuple of actor, object, and action and show that it outperforms baselines approaches that do not capture structured entity-relation information. Ding et al. (2015) leverage the structured event representations of Ding et al. (2014) and learn event embeddings via a neural network designed to output similar embeddings for similar events. The event embeddings enable generalization to unseen events and are used as inputs in deep learning models that forecast individual stock prices. The authors further improve the quality of the event embeddings by using a knowledge graph that provides background knowledge on different entities as well as their relations and show that it leads to more accurate prediction of stock market volatilities (Ding et al. 2016). Our definition of events differs as it captures only the action word(s) of an event, termed event trigger, and we group similar event triggers into an “event class,” which allows us to effectively share information across different occurrences of the same event and address sparsity issues arising from the presence of infrequent or rare events.

3. News Event-Driven Forecasting Model

In this section, we describe our proposed methodology in detail. As mentioned, our method combines event representations extracted from news articles with historical price data using a neural net architecture to forecast future prices. Before we describe each of these components, we formally introduce the problem and the associated notation.

3.1. Problem Definition

Our objective is to build a predictive model that can forecast commodity prices using real-world event occurrences reported in news articles. We assume access to a corpus of news articles indexed by time t , with \mathcal{D}_t denoting the collection of news articles published at time t . The granularity of the time index t depends on the particular prediction task; it can be a day, a week, a month, or a year, and our formulation is agnostic to the actual granularity. The news articles report different real-world events, and we suppose that the reported events come from a fixed but unknown universe of size K . We discuss how to identify these events from a corpus of news articles in Section 3.2. For now, suppose that there exists a function $\phi_t: \mathcal{D}_t \rightarrow [0, 1]^K$ that maps a collection of news articles published at certain time t to a vector $\phi_t(\mathcal{D}_t) = (\phi_{t,1}, \phi_{t,2}, \dots, \phi_{t,K})$ that specifies the “intensity” of each of the K events at time instant t . We call $\phi_t(\mathcal{D}_t)$ the *event embedding* at time t and refer to it as ϕ_t in the remainder of the paper, with the news corpus \mathcal{D}_t being implicit.

Next, let y denote the price time series of any commodity of interest, with y_t representing the price at time t . Our goal is to design a predictive model \mathcal{G} that takes as inputs the event embeddings $\phi_t, \phi_{t-1}, \dots, \phi_{t-\delta+1}$ and the past prices $y_t, y_{t-1}, \dots, y_{t-\delta+1}$ and outputs the price at time $t+1$. Here, $\delta \geq 1$ is a time lag parameter that captures how far in the future a real-world event can influence the price of the commodity. In Section 3.3, we present a neural network-based predictive model for forecasting commodity prices.

Note that in the formulation, we can use any function $\phi_t(\cdot)$ that maps a collection of news articles into a representation of real-world events. Indeed, our evaluation in Section 4 compares the predictive performance of different representation schemes for the mapping ϕ_t .

3.2. Generating Event Embeddings from News Articles

In this section, we formally describe how we generate *event embeddings*, our custom embedding for news articles. The success of our approach critically depends on how well we are able to extract relevant events from the text of a news article. Meaningless “events” could hurt rather than help the prediction accuracy. Because a news article mentions many facts beyond just the events, such as location, entities involved, etc., general-purpose representation schemes, such as LDA, word2vec, or doc2vec, might identify “topics” or themes that are not directly related to the events being reported. These representations do indeed worsen the prediction performance (see Section 4.2).

We address this challenge by constructing event embeddings that exploit the specific structure of news articles. Specifically, our embedding scheme (1) restricts the vocabulary to only consider “event triggers,” and

(2) infers events reported by leveraging the unique discourse structure in news articles. We elaborate on the sequence of steps employed to obtain the event embeddings from the news corpus next.

3.2.1. Automated and Scalable Event Trigger Extraction.

Although the notion of an event is intuitive, algorithmic extraction requires a more precise way to operationalize it. Several approaches have been proposed in the NLP literature for this purpose. Of these, we focus on identifying events using the corresponding event triggers. An event trigger is a word or a phrase appearing in the news article that specifies the occurrence and the type of a specific event (Dodding et al. 2004). For example, from the headline “FIFA Officials Arrested in Corruption Case” of a news article, one may infer that the article is reporting an *arrest* event. The essence of this event can be captured using the word *arrested*, which is the event trigger for this headline. Usually, event triggers are verbs or nouns present in the sentence that describe some notion of “action” or “incident.” Event triggers can occur in various forms—verbs (traders *protest* over FDI in retail), nouns (*burglary* in police station leaves cops red-faced), or a combined phrase (*number of AIDS patients go up* in MP). Ultimately, determining which words or phrases in a sentence are event triggers requires human input. One approach then is to manually label each word of each sentence in a news article as a trigger or a nontrigger. The manual approach, of course, does not scale to the large size of our news corpus (consisting of approximately 1.6 million articles).

To deal with the labeling challenge, we need an algorithm that can mimic human annotators in accomplishing the sequence labeling task, where given an input sentence, the algorithm labels each word in the sentence as a trigger or nontrigger. For this purpose, we propose a supervised learning method using a conditional random field (CRF) model (Lafferty et al. 2001) to automatically extract event triggers from any news article. The CRF model leverages features, such as part of speech and named entity tags, word position, etc., to predict the label (trigger or nontrigger) of each word in an input sentence; we defer the precise details to Online Appendix A.1. To generate data for training the CRF, we conduct an Amazon Mechanical Turk (MTurk) study, described in Online Appendix A.1.1, to annotate a “small” random sample of article headlines with event triggers. The evaluation of our CRF model reveals a high F-1 score of 0.837 in identifying the event triggers, despite being trained with a small amount of labeled data. Our proposed methodology is also versatile as the CRF can easily be trained using different kinds of labeled data in order to extract meaningful representations in alternate contexts.

3.2.2. Clustering Event Triggers to Form Event Classes. To enable information sharing across different occurrences of the same event as well as address data sparsity issues because of rare events, in the second step, we group the event triggers into what we call event classes, where each event class represents a type of event. In particular, we first leverage the trained CRF model to extract event triggers from all news articles in the corpus. Then, we cluster the universe of event triggers obtained using the standard k -means algorithm into K (nonoverlapping) groups, with the distance measure defined as the cosine similarity (Manning et al. 2008) between the word2vec representations of the triggers. The optimal number of clusters was determined according to the “elbow” method (Sugar and James 2003), resulting in $K = 250$ event classes for our corpus.

3.2.3. Generative Event Model Exploiting News Discourse Structure. With the tools described in place, we can identify the numbers of event triggers and event classes mentioned in each news article. But the question remains: what is/are the main event(s) mentioned in the news article? To infer these main events, we next propose a probabilistic generative model for real-world events reported in any news article. Unlike the widely used LDA topic model (Blei et al. 2003) that models each document as a mixture of topics and treats all the words reported in the document equally, our proposed model exploits the unique structure of a news article as outlined.

Although a given news article references multiple events in practice, Choubey et al. (2018) found that there is usually one dominant or *central* event reported in any article. Moreover, the news discourse literature prescribes that the title/headline and the lead paragraph of the article typically summarize the central event being reported; see, for instance, Bell (1991), Van Dijk (2013), and Yarlott et al. (2018). Motivated by these observations, our generative model posits that each news article reports *one* central (or main) event, which is drawn from a finite set of event classes. The event class of a news article, instantiated in the central event, is identified through the event triggers mentioned in the headline and first paragraph of the article. The events referenced in the remainder of the news article are termed *subsidiary events*, where again, each event corresponds to some underlying event class that is closely associated with the central event class. For example, we expect that the event class {“blasts,” “explosion,” “bombing,” ...} is frequently accompanied by the event class {“kill,” “injure,” “die,” ...} in any news article.

We train our generative event model on the entire news corpus using the procedure outlined in Online Appendix A.2. The trained model outputs a probability distribution over the event classes for each news article, ideally assigning the largest probability to the central event reported in the article. To ensure that our model

is identifying meaningful events, we conduct a second MTurk study in which we asked a different set of human annotators to identify the central event reported in news articles. This labeled data set is used to evaluate the generative event model that predicts the most likely central event reported in each news article. Our evaluation, presented in Online Appendix A.2.1, shows that the event model achieves an accuracy of 93.5% and a macro average F-1 score of 0.917, indicating its success in correctly determining the central event in news articles identified by human annotators.

3.2.4. Computing the Embeddings. Finally, the trained generative model is used to infer the most likely central event in each news article. We then aggregate the central event predictions across all articles in \mathcal{D}_t to obtain the event embedding ϕ_t at time t according to (EC.3) in Online Appendix A.2.

3.3. REN: A Neural Network Model for Price Forecasting

We now describe our framework for forecasting commodity prices. Specifically, we introduce the concept of an REN—an RNN-based (Bengio et al. 2017) event-driven predictive model. Neural networks parameterize the input-output relationship using a series of nonlinear transformations (captured by “artificial neurons”) and aim to learn the parameters using large amounts of training data. RNNs are used to capture temporal or sequential dynamics in the data and therefore, are well suited for time-series forecasting (Malhotra et al. 2015, Che et al. 2018). An REN leverages the RNN framework to forecast a time series by capturing nonlinear dependencies on past prices as well as exogenous event information.

Formally, the REN has a sequence of δ hidden states, where the i th hidden state ($1 \leq i \leq \delta$) encodes aggregated event and price information from time $t - \delta + 1$ to $t - \delta + i$. Each hidden state takes as inputs the event embedding and the commodity price at the corresponding time point as well as the output from the previous hidden state. It then applies a state-to-state transition function \mathcal{F} to combine all of its inputs into an output vector $\mathbf{h}_i \in \mathbb{R}^H$ as follows:

$$\mathbf{h}_i = \mathcal{F}(\phi_{t-\delta+i}, y_{t-\delta+i}, \mathbf{h}_{i-1}), \quad (1)$$

where for all $1 \leq i \leq \delta$, $\phi_{t-\delta+i} \in \mathbb{R}^K$ is the event embedding and $y_{t-\delta+i}$ is the commodity price at time $t - \delta + i$. We let $\mathbf{h}_0 = \mathbf{0}$ be the all-zeros vector in \mathbb{R}^H . It is customary to think of each hidden state as comprising H units or neurons, with the values in \mathbf{h}_i representing the “activation” of each unit in hidden state i . The state-to-state transition function \mathcal{F} is chosen as the composition of an element-wise nonlinearity $f: \mathbb{R} \mapsto \mathbb{R}$, such as the logistic sigmoid, hyperbolic tangent function, or rectified

linear unit (Bengio et al. 2017), with an affine transformation of the inputs and the previous hidden state output:

$$h_i = f(\mathbf{W} \cdot [\boldsymbol{\phi}_{t-\delta+i}, y_{t-\delta+i}] + \mathbf{U} \cdot h_{i-1}), \quad (2)$$

where $\mathbf{W} \in \mathbb{R}^{H \times (K+1)}$ is the input-to-hidden state parameter matrix and $\mathbf{U} \in \mathbb{R}^{H \times H}$ is the state-to-state recurrent parameter matrix. Note that the parameters \mathbf{W} and \mathbf{U} are shared across all the hidden states.

The final hidden state outputs the predicted price at time $t + 1$, which we denote as \hat{y}_{t+1} :

$$\hat{y}_{t+1} = \mathbf{v}^\top \mathbf{h}_\delta,$$

where $\mathbf{v} \in \mathbb{R}^H$ represents the hidden state-to-output parameter vector. Note that the predicted price is (implicitly) computed as a function of all the previous event embeddings $\boldsymbol{\phi}_t, \boldsymbol{\phi}_{t-1}, \dots, \boldsymbol{\phi}_{t-\delta+1}$ and all of the past prices $y_t, y_{t-1}, \dots, y_{t-\delta+1}$ because of the sequential dependence in (1).

Summarizing the steps, the end-to-end predictive model for the REN can be written as

$$\hat{y}_{t+1} = \mathcal{G}(\boldsymbol{\phi}_t, \boldsymbol{\phi}_{t-1}, \dots, \boldsymbol{\phi}_{t-\delta+1}, y_t, y_{t-1}, \dots, y_{t-\delta+1}), \quad (3)$$

where \mathcal{G} is the composite function that maps the time series of event embeddings and past prices to the predicted price. The parameters of the REN model are then estimated by minimizing the mean squared error between the actual and predicted prices on historical data, referred to as the *training data set*. The learned parameters are then used to forecast future prices.

Finally, we note that RNNs have been found to suffer from the vanishing gradient problem, which means that the magnitude of the gradient (partial derivative with respect to the model parameters) that is used to update the network during training decreases rapidly as the number of hidden states grows. As a result, they are unable to model longer-term dependencies (Pascanu et al. 2013). For example, some news events might have a “delayed effect” on the price, which the standard RNN model may not be able to capture. The issue was addressed by introducing *memory* to the hidden states. We experimented with both LSTM (Hochreiter and Schmidhuber 1997) and gated recurrent units (Cho et al. 2014) and observed that LSTM units performed better. Therefore, we implement RENs with LSTM units in the hidden states in our numerical evaluation, described next.

4. Evaluation Using a Corpus of News Articles and Price Data

Our data are composed of two different sources over a 15-year period (January 2006 to December 2020)—a corpus of news articles and a set of time-series data for the different commodities. The price data were collected from the Ministry of Agriculture and Farmers Welfare, Government of India’s AgMarknet portal (<https://agmarknet.gov.in>). The Directorate of Marketing and Inspection

under the ministry publishes the minimum, maximum, and modal (the price at which maximum sales were recorded) prices of crops across many different wholesale markets in the country. In this work, we focused on onion, potato, rice, and wheat, which are among the most consumed agricultural products in India. We consider the *daily* modal prices of each commodity and compute the average across all markets in the country between January 1, 2006 and December 31, 2020. This gives us the actual commodity price y_t for each day t . To obtain the event embeddings, we use a corpus of news articles published by the *Times of India* (TOI), which is a leading national daily in India, between 2006 and 2020. We collected all the articles published during this time period from the online archives, which are available at <https://timesofindia.indiatimes.com/archive.cms>. Our corpus consists of 1,684,322 articles from both the national as well as different regional TOI editions. From each article, we extracted the headline, main text, and publishing time.

4.1. Methods Compared

To isolate the benefits of leveraging event information for improved prediction as opposed to employing a more sophisticated (nonlinear) model, we compare the performance of the following different methods for forecasting commodity prices.

1. **Linear models.** We implement the following models that capture linear dependence on past prices and event information.

a. *Naive forecast.* As a simple baseline, we consider the naive forecast that predicts that tomorrow’s price is the same as today’s price: that is,

$$\hat{y}_{t+1} = y_t.$$

b. *Autoregressive (AR) process.* We implement an AR process of order p , which predicts prices according to the following equation:

$$\hat{y}_{t+1} = \alpha_0 + \sum_{\ell=1}^p \alpha_\ell \cdot y_{t+1-\ell}, \quad (4)$$

where $\boldsymbol{\alpha} = (\alpha_0, \alpha_1, \dots, \alpha_p)$ is a vector of parameters that is learned from the data. We experimented with different values for p , and the results were qualitatively similar. We report the results for $p=2$ and use ordinary least squares (OLS) to estimate $\boldsymbol{\alpha}$.

c. *AR process with event embeddings (AR + events).* We update (4) to incorporate event information into the AR(p) model as follows:

$$\hat{y}_{t+1} = \alpha_0 + \sum_{\ell=1}^p \alpha_\ell \cdot y_{t+1-\ell} + \sum_{k=1}^K \omega_k \cdot \left(\sum_{\ell=1}^{\delta} \phi_{t+1-\ell, k} \right), \quad (5)$$

where $\boldsymbol{\omega} = (\omega_1, \omega_2, \dots, \omega_K)$ is a vector of parameters capturing the impact of the different events; recall that

$\phi_{p,k}$ is the event embedding for event $k \in [K]$ for day t' . Similar to RENs, we incorporate the impact of events up to δ days in the past. An alternate specification would be to allow events at different lags to impact the predicted price differently: in other words, estimate a separate coefficient $\omega_{k,\ell}$ for each $k \in [K]$ and $\ell \in [\delta]$. Because $K=250$ in our context, the number of parameters to estimate quickly grows as the event lag δ becomes larger, increasing the risk of overfitting. Consequently, we estimate the parsimonious form of the model outlined in (5). However, we still encountered overfitting in certain scenarios and applied ridge regularization to improve the prediction performance, where the penalty term was chosen using a validation set. Finally, because events can have disparate impacts for different crops, we treat δ as a hyperparameter and rely on crossvalidation to tune its value among $\{1, 7, 14, 21, 28\}$, which correspond to one-day, one-week, two-week, three-week, and four-week lags, respectively.

d. *AR process with alternate embeddings (AR + X).*

We consider three additional models obtained by replacing the event embeddings in Equation (5) with LDA topics (Blei et al. 2003), word2vec embeddings (Mikolov et al. 2013), and doc2vec embeddings (Le and Mikolov 2014). Refer to Online Appendix C for details on how each of these embeddings is computed.

2. **Nonlinear models.** The following models capture nonlinear dependence on past prices and event information.

a. *LSTM model using only past prices.* We consider an LSTM prediction model similar to RENs but that uses only the past prices $y_t, y_{t-1}, \dots, y_{t-\delta+1}$ to forecast the price on day $t+1$. Because training LSTM models is computationally intensive, we set the lag $\delta=7$ days, which provided a reasonable balance between prediction accuracy and model training time. Neural network models, including LSTM, require careful tuning of many hyperparameters as well as appropriate normalization of input features to achieve good performance; we describe these details in Online Appendix B.

b. *REN model.* We fix $\delta=7$ days to analyze the impact of incorporating event embeddings into the pure LSTM model. The training procedure is identical to that used for LSTM.

c. *LSTM model using past prices and alternate embeddings (LSTM + X).* Similar to the AR process, we consider alternate models obtained by replacing the event embeddings in RENs with LDA topics, word2vec, and doc2vec embeddings.

4.2. Results and Discussion

We first generate the event embeddings for the entire news corpus using the procedure described in Section 3.2. Next, we train the different forecasting methods

using data from January 2006 to December 2012. Specifically, we divide this period into overlapping intervals of δ days so that each training instance involves predicting the price y_{t+1} using the event embeddings and historical prices from day $t-\delta$ up to day t . The LSTM and REN models were trained using the TensorFlow (Abadi et al. 2015) library, whereas the AR models were trained using Python's scikit-learn library (Pedregosa et al. 2011). We leverage the trained models to forecast prices for each day in the test period, lasting from January 2013 to December 2020.

Table 1 summarizes the predictive performance by reporting the root mean square error (RMSE) in predicting the one-day-ahead crop price under all competing methods. We make the following observations from the table.

1. *RENs outperform all competing methods.* Our proposed method outperforms the benchmarks for all four crops. In particular, RENs achieve an RMSE reduction of about 7% for onion and potato, 13% for rice, and 5% for wheat over the best benchmark, as reported in Table 1. In fact, the absolute performance of RENs is also impressive; the MAPEs, reported in Table EC.7 in Online Appendix E, for onion, potato, rice, and wheat are 5.2%, 5.3%, 3.4%, and 2.3%, respectively. The high prediction accuracy for onion is particularly promising because onion prices are susceptible to even minor shocks.

2. *Events help improve accuracy—the “right” events even more so.* Both linear and nonlinear models benefit from adding event information, indicating that events contain early warning signals on price changes. In particular, the AR + events method achieves a 5% reduction in RMSE compared with AR for onion price prediction. Similarly, the REN model achieves an average 10% reduction in RMSE over the LSTM benchmark across the four crops. The improvement is also apparent from Figure 3, which plots the daily forecast errors—difference between the predicted and actual prices—for the LSTM and REN models during the test period. It can be seen that the LSTM model exhibits large deviations from zero, indicating significant forecast errors. The signal from event embeddings in the REN model helps to significantly reduce the spread of the daily forecast errors. A similar comparison between AR + events and AR is shown in Figure EC.4 in Online Appendix E.

However, it is important to note that the improvements depend on the specific type of event representation used. Table 1 shows that widely used general-purpose text representations, such as LDA topics (Blei et al. 2003), word embeddings (Mikolov et al. 2013), and document embeddings (Le and Mikolov 2014), can in fact worsen the forecast accuracy because of overfitting. In contrast, event embeddings always improve the performance and provide significant reductions in

Table 1. RMSE in Units of Rupees per Quintal for One-Day-Ahead Price Forecasts During the Test Period (January 2013 to December 2020)

Forecasting methods	Crop			
	Onion	Potato	Rice	Wheat
Linear models				
Naive	188.22	107.00	178.81	65.56
AR	178.05	96.72	161.86	59.11
AR + LDA topics	169.34	96.33	161.60	59.73
AR + word2vec embeddings	171.38	96.50	167.23	62.01
AR + doc2vec embeddings	171.12	98.89	188.50	66.60
AR + events	169.59	96.67	161.78	59.04
Nonlinear models				
LSTM	167.72	91.78	156.99	66.01
LSTM + LDA topics	188.34	107.00	165.10	65.55
LSTM + word2vec embeddings	189.53	107.21	178.94	70.08
LSTM + doc2vec embeddings	189.56	107.08	179.51	66.30
Our method (REN)	155.56	86.08	135.62	55.96

Notes. Incorporating event embeddings helps improve the performance of both linear and nonlinear models, whereas general-purpose embeddings can worsen the forecast accuracy because of overfitting. Nonlinear models outperform linear models, implying that the underlying price fluctuations are complex. Our proposed method, highlighted in bold, reduces the forecast error by about ₹12/quintal for onion (7% lower), ₹6/quintal for potato (7% lower), ₹21/quintal for rice (13% lower), and ₹3/quintal for wheat (5% lower) over the best benchmark.

RMSE compared with alternate embeddings for the nonlinear models: 18% and 20% for potato and rice, respectively.

3. *Nonlinear models can generate more accurate forecasts.* This suggests that the underlying price fluctuations are complex so that linear forecasting models are insufficient in explaining them. For instance, the LSTM model achieves ₹11/quintal lower RMSE than AR for onion price prediction, whereas the event-driven REN model obtains ₹21/quintal reduction in RMSE over the event-based linear model (AR + events) for predicting the price of rice. In fact, from Table 1, we observe that the LSTM model outperforms the AR + events method for three of the crops, which further highlights the benefit of using sophisticated models for predicting commodity prices. At the same time, we observe that the LSTM model performs worse than the AR models and comparable with the naive forecast when predicting wheat prices. This also showcases the need for appropriate regularization of such neural network models when forecasting commodity prices. Including our proposed event embeddings as inputs to the LSTM model provides valuable signals about the factors influencing price changes, serving as a form of regularization.

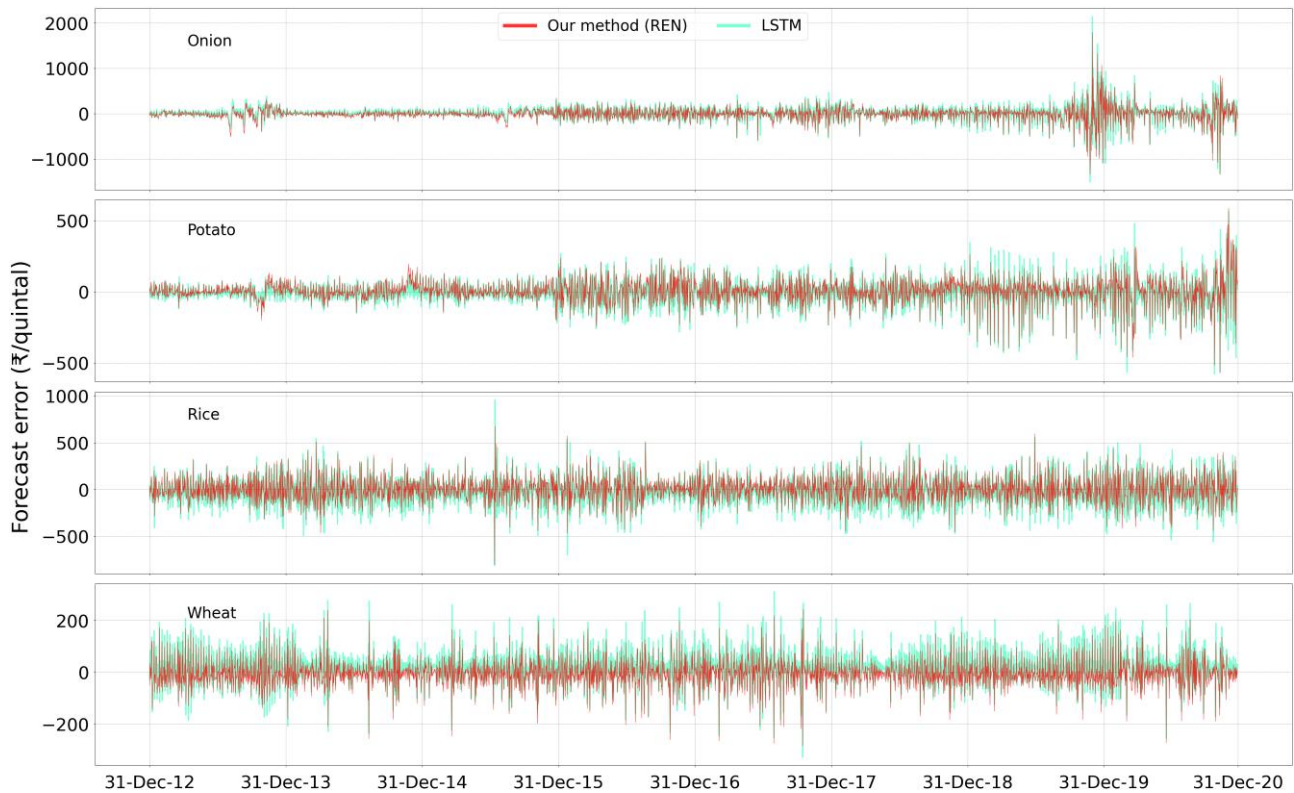
The findings validate our original hypothesis of leveraging real-world event occurrences for predicting commodity price fluctuations, showing that using only historical prices may not provide the most accurate predictions. They also showcase the value of extracting the “right” set of features from text data because our

proposed event embeddings that leverage the special structure of news articles outperform general-purpose embedding techniques.

5. Which Events Are Driving the Price Fluctuations?

The evaluation in the previous section shows the efficacy of using our proposed REN model for forecasting agricultural commodity prices. However, it does not reveal which events (if any) are responsible for the price fluctuations at any particular time. Beyond lending interpretability to our model, which can help build confidence into the model’s predictions (Burkart and Huber 2021), knowledge of such events can assist managers in being better prepared for future occurrences of the same or related events. For instance, if managers know that elections or festivals are typically associated with rising prices, they can influence the procurement strategy before an upcoming election or major holiday. Similarly, it can help investment managers with commodity futures in their portfolios take precautionary actions to hedge against potential losses because of upcoming events. It can also aid governments and policymakers in reducing the negative impacts of impending price shocks by informing key decisions, such as budget allocation and spending, design of subsidies and incentives, etc., prior to such events and ensuring that vulnerable populations have access to affordable food.

To identify the key drivers of price fluctuations at each time, we first compute event importance scores

Figure 3. (Color online) Daily Forecast Error During the Test Period (January 2013 to December 2020)

Notes. Forecast error is computed as the difference between the predicted price and the actual price, so values close to zero are preferred. The LSTM model exhibits large errors during multiple days in the test period. By incorporating event information, our proposed method is able to significantly reduce the spread of the daily forecast errors. The difference between the two methods is visible more clearly in the online color version.

from the trained REN model as described next. Events with the largest importance scores at each time point are then deemed to be the biggest drivers of price fluctuations at that time.

5.1. Computing Event Importance Scores

RNNs leverage the temporal sequence of input values to infer the most likely value of the output variable. As information from the previous inputs flows in the forward direction (via the hidden states), it undergoes several complex transformations, and therefore, it becomes difficult to interpret the *importance* of the input at a specific time point in predicting the output variable. To address this issue, there have been different proposals in the literature to “interpret” a sequential model (such as RNN) and determine the relative importance of a particular input for the prediction outcome. As we move to investigate the important events that drive fluctuations in commodity prices, we surveyed some of these methods and determined the most suitable approach for our case. We specifically concentrated on methods that are applicable to a model post-training, and no information is required from the training phase. Based on applicability of these methods

in our context and the computational complexity, we chose the gradient-based explanation method (Denil et al. 2015, Li et al. 2016) to interpret the relative importance of input features in our REN model.

In the gradient-based explanation method, a *relevance score* is computed to represent the importance of each input variable for predicting the outcome of a single data point. Before we can formally define the relevance score, we need to introduce some notation. For any time t' , let $\mathbf{x}_{t'} = [\boldsymbol{\phi}_{t'}, y_{t'}]$ denote the input vector to the REN model, where recall from Section 3.3 that $\boldsymbol{\phi}_{t'} \in \mathbb{R}^K$ is the K -dimensional event embedding and $y_{t'}$ is the commodity price at time t' . Therefore, the vector $\mathbf{x}_{t'} \in \mathbb{R}^{K+1}$. In our context, a data point is of the form (\mathbf{X}_t, y_t) , where y_t is the commodity price at time t and $\mathbf{X}_t = \{\mathbf{x}_{t-1}, \dots, \mathbf{x}_{t-1-\delta}\}$ is the set of input vectors in the interval $\mathcal{T}_{t,\delta} := [t-1-\delta, t-1]$. An input variable thus corresponds to any of the features $x_{t',i}$, where $t' \in \mathcal{T}_{t,\delta}$ and $1 \leq i \leq K+1$. For any input variable $x_{t',i}$, we denote by $R_{t',i}^{(t)}$ its relevance score, defined as $R_{t',i}^{(t)} = \left| \frac{\partial \mathcal{G}(\mathbf{X}_t)}{\partial x_{t',i}} \right| \times |x_{t',i} - \mathbb{E}[x_{t',i}]|$, where $\mathbb{E}[x_{t',i}]$ denotes the expected value of feature $x_{t',i}$ and \mathcal{G} is the composite mapping from (3) that computes the predicted price. Note that the relevance score is nonnegative and is

defined as the product of two terms: the derivative that captures how much the predicted price changes with a small change in the input feature—the larger the change in the predicted price, the more relevant it is—and the relative intensity of the input feature measured by its deviation from the average value.

Given this, the relevance score $R_i^{(t)}$ of feature i for the data point (X_t, y_t) is computed by summing up the relevance scores for that feature across all the input units to the REN; that is, we compute $R_i^{(t)} := \sum_{t' \in \mathcal{T}_{t, \delta}} R_{t', i}^{(t)}$. Then, for each $1 \leq i \leq K$, the score $R_i^{(t)}$ captures the relevance or importance of event i in predicting the price at time t . In a similar fashion, $R_{K+1}^{(t)}$ captures the importance of the historical price. Finally, for each feature i , we compute its importance in any month M , say $R_i^{(M)}$, as the average of the daily importance scores, $R_i^{(M)} = \sum_{t \in M} R_i^{(t)} / \text{num_days}(M)$, where $\text{num_days}(M)$ is the number of days in month M . We choose monthly (instead of daily) intervals to avoid focusing on noisy local fluctuations in the prices.

5.2. Observations and Insights

In Figure 4, we present a heat map for each crop showing the relative importance of different events in predicting price fluctuations in each month during a subset of the test period (January 2013 to December 2015); we only show the subset of events that were identified to be among the five most important events in at least one month. Each event is manually labeled with a single word that best describes the event; see Online Appendix D for the collection of event triggers associated with each label. For onion, potato, wheat, and rice, our analysis identifies 11, 12, 12, and 11 unique events, respectively, as the key drivers of price fluctuations in the three-year period. We make the following observations.

1. Events driving price fluctuations for a particular crop vary over time. We find that for all four crops, the importance of any particular event varies across the test period, and therefore, different events drive fluctuations in crop prices at different times. For instance, although the *natural* event is found to be among the most important factors behind rice price fluctuations (TOI 2013a, 2014a; Economist 2015) during the summer months of 2013–2015, the event *scam* was the main driver of rice prices during late 2013 and early 2014 (TOI 2013b, 2014b, c).

2. Different events drive price fluctuations for different crops during the same time period. Our analysis finds that even during the same time period, the events driving price fluctuations differ across different crops. For instance, in contrast to the *natural* event for rice mentioned earlier, wheat prices were mostly influenced by the *policy* event during the summers of 2013 and 2015 (TOI 2013c, 2015a; Business Standard 2015).

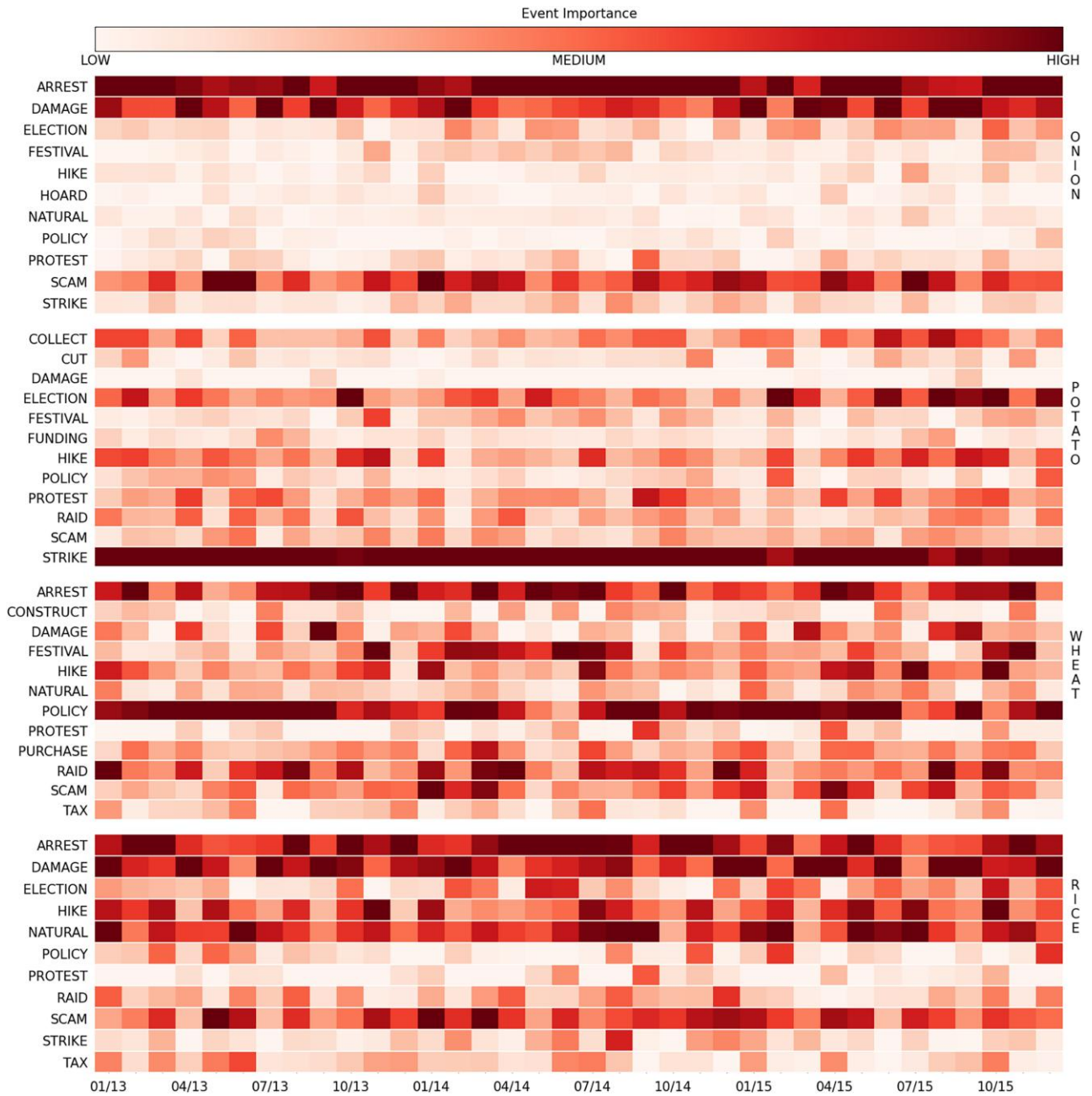
Similarly, during November 2013, wheat price changes were attributed to the *festivals* prevalent during that period. On the other hand, the event *hike*, which refers to a hike in prices of other commodities, such as fuel and electricity, or transportation-related factors (e.g., fare hike or hike in highway tolls), was deemed to be the most important driver of rice prices in the same month. Although we did not find any direct evidence of such hikes influencing the crop prices, there are news articles reporting fuel and transportation price hikes in the months leading up to November 2013 (Bagri 2013; TOI 2013d, e), which may have impacted the supply chain by increasing transportation costs, thereby driving prices up.

3. Our analysis identifies well-known factors driving commodity prices. Past studies have shown that oil prices, climatic conditions, and government policies (e.g., subsidies) lead to major changes in food prices (Amadeo 2021, Bogmans et al. 2021). These correspond to the event classes *hike*, *natural*, and *policy*, respectively. Existing studies have also identified the availability of resources, such as water, as one of the key factors in ensuring long-term food security (Grafton et al. 2015). Although our coarse-grained event classes do not identify specific events, the class labeled as *natural* (which encompasses events such as poor monsoon/rainfall, floods, drought, etc.) broadly explains the impact of water availability on commodity prices. Similarly, the event classes *strike*, *protest*, and *damage* highlight social unrest as a plausible cause for increased commodity prices, a finding corroborated by several prior studies (Bellemare 2015, Raleigh et al. 2015, Li et al. 2022).

The findings show that our proposed REN model is able to associate real-world event occurrences with fluctuations in the commodity prices and is not merely detecting spurious correlations. However, it is worth noting that our method provides only a coarse-grained picture by identifying a candidate set of factors that are most predictive of the crop prices, after which a more fine-grained analysis can be conducted to understand the precise impact of different events on the prices. For instance, one could examine news articles published in the corresponding time interval and look for mentions of the events identified by our analysis. We performed such an exercise for a subset of events identified in Figure 4, which revealed some interesting insights.

- *Natural* events were found to be among the most important factors driving price fluctuations for rice and wheat. Examples include floods (India Today 2015) and excessive or poor rainfall (Dutta 2014, Bera 2015, TOI 2015b, WorldGrain 2015). Such events can damage the crop harvest and adversely impact supply and production, subsequently driving the prices up.

Figure 4. (Color online) Variation in the Importance of Different Events in Driving the Underlying Price Fluctuations Between January 2013 and December 2015



Notes. We compute the importance of an event using a gradient-based method that measures the sensitivity of the predicted price to its intensity. The importance scores are normalized to lie between zero (low) and one (high). For each crop, the scores are displayed in a heat map, with the row corresponding to an event and the column corresponding to a month. Our analysis identifies events known to impact price fluctuations, such as elections, festivals, government policies, and natural calamities. For each crop, we see that the importance of any event varies over time, indicating that different events are driving the price fluctuations during different times. Moreover, even during the same time period, different events may be responsible for driving the price fluctuations of different crops.

- Similarly, food prices were uniformly found to be impacted by various policy measures taken by the government—denoted by the event *policy*—to better match supply and demand. We notice that such policies seem to take effect when the prices have been

steadily increasing for a few months and the government is under pressure to take action to alleviate the issue. For instance, potato prices soared during the first half of 2014 and 2015. As a result, the government took some steps to control the price rise, such as curbing

exports and providing subsidies to lower transport costs (KNN 2014; TOI 2014d, 2015c; Acharya 2015; Economic Times 2015; Frontline 2015).

- We also found some novel events responsible for the price fluctuations in many months, such as *scam*. Upon further investigation, we found that numerous illicit practices have caused food prices to increase (TOI 2013b, TheHindu 2015, India Today 2015a). A similar practice is *hoarding* of onion stock by traders, which is often highlighted in the media as a means to artificially raise the prices (Mukherjee 2014, TOI 2014e, The Hindu Businessline 2020). In fact, even legally hoarding can affect the prices (Jagannathan 2014). On the contrary, the event *raid* counters this effect and is found to be one of the important factors driving the prices, backed by several news reports of law enforcement agencies raiding storage units across India to take action against hoarders (TOI 2015d, e, f). Another example is the event *purchase* identified for wheat, which refers to various steps taken for more systematic procurement of the produce (Economic Times 2014, Nibber 2021).

- Finally, some events were consistently found to be key drivers throughout the duration considered, such as *arrest* for onion and *strike* for potato. Transportation strikes directly impact food prices, including for potato (India Today 2015b; TOI 2015g, 2018), and are likely to be an important driver. On the other hand, an event such as *arrest* may not seem to have a direct impact, but because of its ubiquitous nature, it can act as a confounding or mediating driver of food prices. This suggests the potential of studying the cascading impact of persistent events in driving commodity prices in future work.

6. Conclusions, Limitations, and Future Directions

Based on the premise that prices of many commodities are sensitive to real-world events, this paper presents a data-driven framework to extract events reported in news articles and incorporate them as exogenous features to improve the forecasting accuracy. We introduce the notion of *RENs*, a neural network-based framework for building event-driven nonlinear predictive models for forecasting commodity prices, and evaluate their performance in forecasting one-day-ahead prices of four staple agricultural commodities in India. Our method achieved up to 16% improvement in forecast accuracy when compared with baseline methods that employ linear predictive models and 20% improvement compared with using general-purpose event representations. These results show the benefits of using nonlinear models as well as the “right” event representation scheme for predicting commodity prices. Our framework also provides interpretability to the generated forecasts by identifying the most important events driving the price changes during each time period.

Accurate commodity price forecasts can benefit individuals, businesses, governments, and the society as a whole (Assefa et al. 2015). Agribusinesses can use precise forecasts to make informed decisions about when and how much to procure, enabling them to optimize inventory and storage costs effectively. The same applies to other players in the food industry, including restaurants and hotels, that can proactively manage their costs, negotiate better contracts with suppliers, and make strategic decisions to reduce financial risk associated with the volatile prices. For farmers, access to reliable price predictions can serve as an invaluable tool to time their harvests effectively, choosing to sell when prices are high to maximize profits. Additionally, if market-level forecasts are available, farmers can further boost their income by strategically choosing markets to sell their current harvest by factoring in transportation and other expenses. It can also empower them with better negotiating power with buyers, processors, and distributors and help them avoid selling their produce at disadvantageous prices. Accurate forecasts can also help in reducing wastage and ensuring more sustainable farming practices. Finally, forecasts of food prices can inform consumers about market trends and enable them to budget for food-related expenses effectively.

Our paper adds to the growing literature on leveraging news articles to improve the forecasting of other important indicators, such as infectious disease spread (Bhatia et al. 2021, Kim and Ahn 2021); macroeconomic variables, like GDP, unemployment, industrial production, inflation, etc. (Feuerriegel and Gordon 2019, Kalamara et al. 2022, Barbaglia et al. 2023); and food (in-)security (Ba et al. 2022, Balashankar et al. 2023). Moreover, our work can lay the foundation for a new research area on developing novel methods for leveraging other kinds of unstructured data, such as social media posts, images, videos, etc., to improve operational decisions.

Although our proposed method is largely flexible and can be applied across various domains, it does suffer from a few limitations. Currently, we use news data from a single albeit reputed source, which might raise questions about the accuracy or potential systemic bias of the information reported in the news articles. Aggregating news across multiple sources can help mitigate such concerns. Further, our framework only leverages English-language articles, possibly missing some local or regional events. By including texts from local language publications, the breadth of reported events could be enhanced. Finally, our approach is best suited for forecasting variables that are *directly* impacted by real-world events. For instance, we would expect it to perform better in predicting the prevalence of seasonal and nonchronic diseases (e.g., flu, malaria, etc.), of which news media often has extensive coverage, thus providing early warning signs of an outbreak. This is in contrast to chronic diseases, like cancer or diabetes,

where the number of cases tends to be more uniform over time.

In addition to addressing the limitations, there are numerous other avenues for future work. First, we note that in general, our method only finds *associations* between events and price fluctuations, not necessarily *causal* links. These associations improve recall by allowing us to construct a small set of candidate events from a large corpus of possible events. Subsequent studies can be performed to find causal links between the candidate events and price fluctuations using methods similar to those of Kang et al. (2017) and Balashankar et al. (2019). Second, our current embedding scheme does not capture the sequence of the event triggers present in an article. Potential improvements include incorporating the context—event triggers surrounding a particular trigger—and modeling the dependence between specific triggers using an attention mechanism, which allows the model to focus on specific parts of the input (Vaswani et al. 2017). Third, in this paper, we considered each commodity in isolation and trained a separate predictive model. A promising avenue for future research is leveraging multitask learning techniques to capture dependencies between different commodities to further improve the prediction performance. Fourth, there is an opportunity to combine our event-driven forecasting approach with other forecasting models. For instance, experts can leverage the knowledge of events impacting commodity prices to improve their own forecasts, which are often included in prediction models (e.g., Cortazar et al. 2019). In particular, incorporating events can help reduce the volatility that is often observed in expert forecasts. Similarly, although we focused on single-regime forecasting models in this work, our event-driven approach can also be applied in conjunction with regime-switching models that are employed when dealing with external shocks. Finally, our current event model ignores the role of specific entities, such as people, organizations, political parties, etc., in driving the commodity prices; incorporating them could lead to further improvements in the forecasting accuracy.

Acknowledgments

The authors thank the department editor, the associate editor, and two anonymous referees for their valuable comments that improved the manuscript. The authors are also grateful to Profs. Rohit Deo and Andrew Wu for insightful discussions and feedback that helped shape the final manuscript. The work also benefited from enriching conversations with participants of the 2023 *Manufacturing & Service Operations Management* Interface of Finance, Operations and Risk Management Special Interest Group Meeting. L. Subramanian is a co-founder of Entrupy Inc, Gaius Networks Inc, and has served as a consultant for the World Bank and the Governance Lab. S. Chakraborty, S. Jagabathula, L. Subramanian are co-founders of Velai Inc.

References

- Abadi M, Agarwal A, Barham P (2015) TensorFlow: Large-scale machine learning on heterogeneous systems. Accessed February 24, 2024, <https://www.tensorflow.org/>.
- Acharya N (2015) West Bengal plans freight subsidy for potato export. Accessed February 24, 2024, https://www.business-standard.com/article/markets/west-bengal-plans-freight-subsidy-for-potato-export-115031400437_1.html.
- Agarwal K (2019) Will anger over low potato prices harm BJP in UP? Accessed February 24, 2024, <https://thewire.in/agriculture/bjp-up-elections-2019-potato-prices>.
- Amadeo K (2021) Why food prices are rising, recent trends, and 2021 forecast. Accessed February 24, 2024, <https://www.thebalance.com/why-are-food-prices-rising-causes-of-food-price-inflation-3306099>.
- Assefa TT, Meuwissen MP, Oude Lansink AG (2015) Price volatility transmission in food supply chains: A literature review. *Agribusiness* 31(1):3–13.
- Atefeh F, Khreich W (2015) A survey of techniques for event detection in Twitter. *Computational Intelligence* 31(1):132–164.
- Ba CT, Choquet C, Interdonato R, Roche M (2022) Explaining food security warning signals with YouTube transcriptions and local news articles. *Proc. 2022 ACM Conf. Inform. Tech. Social Good* (Association for Computing Machinery, New York), 315–322.
- Bagozzi BE, Schrodt PA (2012) The dimensionality of political news reports. *Proc. Eur. Political Sci. Assoc. Meetings, Berlin*.
- Bagri NT (2013) Food and fuel prices push Indian inflation higher. *The New York Times* (November 14), <https://www.nytimes.com/2013/11/15/business/food-and-fuel-prices-push-indian-inflation-higher.html>.
- Balashankar A, Subramanian L, Fraiberger SP (2023) Predicting food crises using news streams. *Sci. Adv.* 9(9):eabm3449.
- Balashankar A, Chakraborty S, Fraiberger S, Subramanian L (2019) Identifying predictive causal factors from news streams. *Proc. 2019 Conf. Empirical Methods Natural Language Processing 9th Internat. Joint Conf. Natural Language Processing (EMNLP-IJCNLP)* (Association for Computational Linguistics, Kerrville, TX), 2338–2348.
- Ban GY, Rudin C (2019) The big data newsvendor: Practical insights from machine learning. *Oper. Res.* 67(1):90–108.
- Barbaglia L, Consoli S, Manzan S (2023) Forecasting with economic news. *J. Bus. Econom. Statist.* 41(3):708–719.
- Bell A (1991) *The Language of News Media* (Blackwell, Oxford, UK).
- Bellemare MF (2015) Rising food prices, food price volatility, and social unrest. *Amer. J. Agricultural Econom.* 97(1):1–21.
- Bengio Y, Goodfellow I, Courville A (2017) *Deep Learning*, vol. 1 (MIT Press, Cambridge, MA).
- Bera S (2015) Wheat output may decline by up to 5. Accessed February 24, 2024, <https://www.livemint.com/Politics/XoPy0NSXUgfbYctqESAQAM/Indias-2015-wheat-output-could-fall-by-45-says-farmmini.html>.
- Berling P, Martínez-de Albéniz V (2011) Optimal inventory policies when purchase price and demand are stochastic. *Oper. Res.* 59(1):109–124.
- Bhatia S, Lassmann B, Cohn E, Desai AN, Carrion M, Kraemer MU, Herringer M, et al. (2021) Using digital surveillance tools for near real-time mapping of the risk of infectious disease spread. *NPJ Digital Medicine* 4(1):73.
- Birthal P, Negi A, Joshi P (2019) Understanding causes of volatility in onion prices in India. *J. Agribusiness Developing Emerging Econom.* 9(3):255–275.
- Blei DM, Ng AY, Jordan MI (2003) Latent Dirichlet allocation. *J. Machine Learn. Res.* 3:993–1022.
- Bogmans C, Pescatori A, Prifti E (2021) Four facts about soaring consumer food prices. *IMF Blog* (June 24), <https://blogs.imf.org/2021/06/24/four-facts-about-soaring-consumer-food-prices/>.

- Boone T, Ganeshan R, Hicks RL, Sanders NR (2018) Can Google trends improve your sales forecast? *Production Oper. Management* 27(10):1770–1774.
- Brockwell PJ, Davis RA (2002) *Introduction to Time Series and Forecasting*, 2nd ed. (Springer, Berlin).
- Burkart N, Huber MF (2021) A survey on the explainability of supervised machine learning. *J. Artificial Intelligence Res.* 70:245–317.
- Business Standard (2015) Maggi ban: Flour millers recall maida supplied to nestle. Accessed February 24, 2024, https://www.business-standard.com/article/companies/maggi-ban-flour-millers-recall-maida-supplied-to-nestle-115061600725_1.html.
- Bybee L, Kelly BT, Manela A, Xiu D (2021) Business news and business cycles. NBER Working Paper No. 29344, National Bureau of Economic Research, Cambridge, MA.
- Carbonneau R, Laframboise K, Vahidov R (2008) Application of machine learning techniques for supply chain demand forecasting. *Eur. J. Oper. Res.* 184(3):1140–1154.
- Caselli T, Vossen P (2017) The event storyline corpus: A new benchmark for causal and temporal relation extraction. *Proc. Events Stories News Workshop* (Association for Computational Linguistics, Kerrville, TX), 77–86.
- Chakraborty S, Subramanian L (2016) Extracting signals from news streams for disease outbreak prediction. *2016 IEEE Global Conf. Signal Inform. Processing (GlobalSIP)* (IEEE, Piscataway, NJ), 1300–1304.
- Chakraborty S, Venkataraman A, Jagabathula S, Subramanian L (2016) Predicting socio-economic indicators using news events. *KDD 2016 Proc. 22nd ACM SIGKDD Internat. Conf. Knowledge Discovery Data Mining* (Association for Computing Machinery, New York), 1455–1464.
- Chambers N (2013) Event schema induction with a probabilistic entity-driven model. *Proc. 2013 Conf. Empirical Methods Natural Language Processing* (Association for Computational Linguistics, Kerrville, TX), 1797–1807.
- Chambers N, Jurafsky D (2011) Template-based information extraction without the templates. *Proc. 49th Annual Meeting Assoc. Comput. Linguistics Human Language Tech. Vol. 1* (Association for Computing Machinery, New York), 976–986.
- Che Z, Purushotham S, Cho K, Sontag D, Liu Y (2018) Recurrent neural networks for multivariate time series with missing values. *Sci. Rep.* 8(1):1–12.
- Cheung JCK, Poon H, Vanderwende L (2013) Probabilistic frame induction. *Proc. 2013 Conf. North Amer. Chapter Assoc. Comput. Linguistics Human Language Tech.* (Association for Computational Linguistics, Kerrville, TX), 837–846.
- Cho K, van Merriënboer B, Gulcehre C, Bahdanau D, Bougares F, Schwenk H, Bengio Y (2014) Learning phrase representations using RNN encoder–decoder for statistical machine translation. Moschitti A, Pang B, Daelemans W, eds. *Proc. 2014 Conf. Empirical Methods Natural Language Processing (EMNLP)* (Association for Computational Linguistics, Kerrville, TX), 1724–1734.
- Choubey PK, Raju K, Huang R (2018) Identifying the most dominant event in a news article by mining event coreference relations. *Proc. 2018 Conf. North Amer. Chapter Assoc. Computat. Linguistics Human Language Tech.*, vol. 2 (Association for Computational Linguistics, Kerrville, TX), 340–345.
- CNBC-TV18 (2019) Onion prices, the slayer of many governments, on the boil again as elections approach. Accessed February 24, 2024, <https://www.cnbctv18.com/agriculture/onion-prices-the-slayer-of-many-governments-on-the-boil-again-as-elections-approach-4412221.htm>.
- Cortazar G, Millard C, Ortega H, Schwartz ES (2019) Commodity price forecasts, futures prices, and pricing models. *Management Sci.* 65(9):4141–4155.
- Cui R, Gallino S, Moreno A, Zhang DJ (2018) The operational value of social media information. *Production Oper. Management* 27(10):1749–1769.
- De Janvry A, Sadoulet E (2009) The impact of rising food prices on household welfare in India. UC Berkeley: Institute for Research on Labor and Employment. Retrieved February 24, <https://escholarship.org/uc/item/7xj9n1qq>.
- Denil M, Demiraj A, De Freitas N (2015) Extraction of salient sentences from labelled documents. Preprint, submitted December 21, <https://arxiv.org/abs/1412.6815>.
- Deuskar N (2020) Delhi election: How onion prices led to Swaraj gov't's defeat in 1998 — The last time BJP held reigns in national capital. <https://www.moneycontrol.com/news/politics/delhi-election-how-onion-prices-led-to-swaraj-govts-defeat-in-1998-the-last-time-bjp-held-reignsin-national-capital-4897921.html>.
- Dev SM (2011) Rising food crisis and financial crisis in India: Impact on women and children and ways of tackling the problem, Indira Gandhi Institute of Development Research, Mumbai Working Papers 2011-003, Indira Gandhi Institute of Development Research, Mumbai, India.
- Digital Green (2020) Personal communication, August 2020. Accessed February 24, 2024, <https://www.digitalgreen.org/>.
- Ding X, Zhang Y, Liu T, Duan J (2014) Using structured events to predict stock price movement: An empirical investigation. *Proc. 2014 Conf. Empirical Methods Natural Language Processing (EMNLP)* (Association for Computational Linguistics, Kerrville, TX), 1415–1425.
- Ding X, Zhang Y, Liu T, Duan J (2015) Deep learning for event-driven stock prediction. *Proc. 24th Internat. Conf. Artificial Intelligence (AAAI Press, Palo Alto, CA)*, 2327–2333.
- Ding X, Zhang Y, Liu T, Duan J (2016) Knowledge-driven event embedding for stock prediction. Yuji M, Rashmi P, eds. 2016, *26th Int. Conf. Comput. Linguistics: Tech. Papers* (The COLING 2016 Organizing Committee, Osaka, Japan), 2133–2142.
- Doddington G, Mitchell A, Przybocki M, Ramshaw L, Strassel S, Weischedel R (2004) The automatic content extraction (ACE) program tasks, data, and evaluation. *Proc. Fourth Internat. Conf. Language Resources Evaluation (LREC'04)* (European Language Resources Association (ELRA), Lisbon, Portugal).
- Dutta R (2014) Late monsoon starts Indian farmer's 'journey to hell'. Accessed February 24, 2024, <https://www.reuters.com/article/us-india-monsoon-farmers/late-monsoon-starts-indian-farmers-journey-tohell-idUKKBN0FP03720140720>.
- Economic Times* (2014) ITC, Cargill and other private companies increase wheat procurement. (May 6), <https://economictimes.indiatimes.com/markets/commodities/itc-cargill-and-other-private-companies-increase-wheat-procurement/articleshow/34709003.cms>.
- Economic Times* (2015) West Bengal's potato farmers in trouble as other states grow their own. (March 20), <https://economictimes.indiatimes.com/markets/commodities/west-bengals-potato-farmers-in-trouble-as-other-states-grow-their-own/articleshow/46628935.cms>.
- Economist* (2015) Of rainfall and price rises. (June 25), <https://www.economist.com/finance-and-economics/2015/06/25/of-rainfall-and-price-rises>.
- Feuerriegel S, Gordon J (2019) News-based forecasts of macroeconomic indicators: A semantic path model for interpretable predictions. *Eur. J. Oper. Res.* 272(1):162–175.
- Frontline* (2015) Harvest of tragedy. (April 29), <https://frontline.thehindu.com/the-nation/harvest-of-tragedy/article7150453.ece>.
- Gidofalvi G (2001) Using news articles to predict stock price movements. Technical report, Department of Computer Science and Engineering, University of California, San Diego.
- Goel A, Gutierrez GJ (2011) Multitechelon procurement and distribution policies for traded commodities. *Management Sci.* 57(12):2228–2244.
- Grafton RQ, Daugbjerg C, Qureshi ME (2015) Toward food security by 2050. *Food Security* 7(2):179–183.
- Gro Intelligence (2016) When onion prices bring Indians to tears. Accessed February 24, 2024, <https://gro-intelligence.com/insights/articles/onion-prices-in-india>.

- Hagenau M, Liebmann M, Neumann D (2013) Automated news reading: Stock price prediction based on financial news using context-capturing features. *Decision Support Systems* 55(3):685–697.
- Haksöz C, Seshadri S (2007) Supply chain operations in the presence of a spot market: A review with discussion. *J. Oper. Res. Soc.* 58(11):1412–1429.
- Headey D, Fan S (2008) Anatomy of a crisis: The causes and consequences of surging food prices. *Agricultural Econom.* 39(s1):375–391.
- Hochreiter S, Schmidhuber J (1997) Long short-term memory. *Neural Comput.* 9(8):1735–1780.
- Hogenboom F, Frasincar F, Kaymak U, De Jong F (2011) An overview of event extraction from text. Erp M, van Hage WR, van Hollink L, Jameson A, Troncy R, eds. *Proc. Detection, Representation, Exploitation Events Semantic Web (DeRiVE 2011), Workshop Conjunction 10th Internat. Semantic Web Conf. 2011 (ISWC 2011)* (CEUR-WS.org, Aachen, Germany), 48–57.
- Hu Z, Liu W, Bian J, Liu X, Liu TY (2018) Listening to chaotic whispers: A deep learning framework for news-oriented stock trend prediction. *Proc. Eleventh ACM Internat. Conf. Web Search Data Mining* (ACM, New York), 261–269.
- India.com (2019) As onion price rise to all-time high, people resort to robbery of kitchen staple in UP, West Bengal. (December 10), <https://www.india.com/news/india/as-onion-price-rise-to-all-time-high-peoplereport-to-robbery-of-kitchen-staple-in-up-west-bengal-3873597/>.
- India Today (2015) Farmers face heavy losses due to untimely rainfall. (March 3), <https://www.indiatoday.in/mail-today/story/farmers-face-heavy-losses-due-to-untimely-rainfall-242737-2015-02-13>.
- India Today (2015a) The great Indian dal scam revealed. (November 24), <https://www.indiatoday.in/india/story/the-greatindian-dal-scam-revealed-273540-2015-11-23>.
- India Today (2015b) Telangana: Trucks go off road, goods movement hit. (June 24), <https://www.indiatoday.in/india/story/telangana-trucks-go-off-the-road-indefinite-strike-259348-2015-06-24>.
- International Trade Administration (2022) Food and agriculture value chain. Accessed February 24, 2024, <https://www.trade.gov/country-commercial-guides/india-food-and-agriculture-value-chain>.
- Jagannathan R (2014) How to turn FCI from food hoarder to food manager – and save tonnes of money. (December 21), <https://www.firstpost.com/business/economy/how-to-turn-fci-from-food-hoarder-to-food-manager-and-save-tonnes-of-money-1985065.html>.
- Kalamara E, Turrell A, Redl C, Kapetanios G, Kapadia S (2022) Making text count: Economic forecasting using newspaper text. *J. Appl. Econometrics* 37(5):896–919.
- Kang D, Gangal V, Lu A, Chen Z, Hovy E (2017) Detecting and explaining causes from text for a time series event. Palmer M, Hwa R, Riedel S, eds. *Proc. 2017 Conf. Empirical Methods Natural Language Processing* (Association for Computational Linguistics Copenhagen, Denmark), 2758–2767.
- Kawamae N (2011) Trend analysis model: Trend consists of temporal words, topics, and timestamps. *Proc. Fourth ACM Internat. Conf. Web Search Data Mining* (ACM, New York), 317–326.
- Kim J, Ahn I (2021) Infectious disease outbreak prediction using media articles with machine learning models. *Sci. Rep.* 11(1):1–13.
- KNN (2014) Curbs on potato shipments; minimum export price fixed at USD 450. (June 2), <https://knnindia.co.in/news/newsdetails/economy/curbs-on-potato-shipments-minimum-export-price-fixed-at-usd-450>.
- Kurawarwala AA, Matsuo H (1996) Forecasting and inventory management of short life-cycle products. *Oper. Res.* 44(1):131–150.
- Lafferty JD, McCallum A, Pereira FC (2001) Conditional random fields: Probabilistic models for segmenting and labeling sequence data. *Proc. Eighteenth Internat. Conf. Machine Learn.* (Morgan Kaufmann Publishers Inc., San Francisco, CA), 282–289.
- Larsen VH, Thorsrud LA (2019) The value of news for economic developments. *J. Econometrics* 210(1):203–218.
- Le Q, Mikolov T (2014) Distributed representations of sentences and documents. Xing EP, Jebara T, eds. *Internat. Conf. Machine Learn.*, Proceedings of Machine Learning Research, vol. 22 (PMLR, New York), 1188–1196.
- Levi R, Singhvi S, Zheng Y (2022) Artificial shortage in agricultural supply chains. *Manufacturing Service Oper. Management* 24(2):746–765.
- Li J, Chen X, Hovy E, Jurafsky D (2016) Visualizing and understanding neural models in NLP. Knight K, Nenkova A, Rambow O, eds. *Proc. 2016 Conf. North American Chapter Assoc. Comput. Linguistics Human Language Tech* (Association for Computational Linguistics, San Diego), 681–691.
- Li X, Xie H, Chen L, Wang J, Deng X (2014) News impact on stock price return via sentiment analysis. *Knowledge-Based Systems* 69: 14–23.
- Li XY, Li X, Fan Z, Mi L, Kandakji T, Song Z, Li D, Song XP (2022) Civil war hinders crop production and threatens food security in Syria. *Nature Food* 3(1):38–46.
- Liu Q, Cheng X, Su S, Zhu S (2018) Hierarchical complementary attention network for predicting stock price movements with news. *Proc. 27th ACM Internat. Conf. Inform. Knowledge Management* (ACM, New York), 1603–1606.
- Madaan L, Sharma A, Khandelwal P, Goel S, Singla P, Seth A (2019) Price forecasting & anomaly detection for agricultural commodities in India. *Proc. 2nd ACM SIGCAS Conf. Comput. Sustainable Soc.* (ACM, New York), 52–64.
- Mahajan S, Sousa-Poza A, Datta K (2015) Differential effects of rising food prices on Indian households differing in income. *Food Security* 7(5):1043–1053.
- Malhotra P, Vig L, Shroff G, Agarwal P (2015) Long short term memory networks for anomaly detection in time series. *ESANN 2015 Proc., Eur. Sympos. Artificial Neural Networks, Comput. Intelligence Machine Learning, Bruges, Belgium*, 22–24.
- Mandl C (2019) Optimal procurement and inventory control in volatile commodity markets. PhD thesis, Technische Universität, Munich, Germany.
- Mandl C, Minner S (2023) Data-driven optimization for commodity procurement under price uncertainty. *Manufacturing Service Oper. Management* 25(2):371–390.
- Mandl C, Nadarajah S, Minner S, Gavirneni S (2022) Data-driven storage operations: Cross-commodity backtest and structured policies. *Production Oper. Management* 31(6):2438–2456.
- Manning CD, Raghavan P, Schütze H (2008) *Introduction to Information Retrieval* (Cambridge University Press, Cambridge, UK).
- McLain S (2013) India's onion prices play pivotal role in election. Accessed February 24, 2024, <https://www.wsj.com/articles/SB10001424127887324906304579038671391636800>.
- Mikolov T, Sutskever I, Chen K, Corrado GS, Dean J (2013) Distributed representations of words and phrases and their compositionality. Burges CJ, Bottou L, Welling M, Ghahramani Z, Weinberger KQ, eds. *Adv. Neural Inform. Processing Systems*, Distributed Representations of Words and Phrases and their Compositionality, vol. 26 (Curran Associates, Inc., Red Hook, NY), 3111–3119.
- Ming F, Wong F, Liu Z, Chiang M (2014) Stock market prediction from WSJ: Text mining via sparse matrix factorization. *2014 IEEE Internat. Conf. Data Mining* (IEEE, Piscataway, NJ), 430–439.
- Mittal S, Hariharan VK, Subash S (2018) Price volatility trends and price transmission for major staples in India. *Agricultural Econom. Res. Rev.* 31(1):65–74.
- Mukherjee S (2014) Middlemen responsible for food inflation: Ram Vilas Paswan. Accessed February 24, 2024, https://www.business-standard.com/article/economy-policy/middlemen-responsible-for-food-inflation-ram-vilas-paswan-114081100014_1.html.
- Nassirtoussi AK, Aghabozorgi S, Wah TY, Ngo DCL (2014) Text mining for market prediction: A systematic review. *Expert Systems Appl.* 41(16):7653–7670.

- New York Times* (2013) Rising onion prices tempt highway robbers in India. Accessed February 24, 2024, <https://india.blogs.nytimes.com/2013/08/21/rising-onion-prices-tempt-highway-robbers-in-india/>.
- Nibber GS (2021) Punjab to put a cap on per acre yield of paddy. *Hindustan Times* (October 6), <https://www.hindustantimes.com/cities/chandigarh-news/punjab-to-put-a-cap-on-per-acre-yield-of-paddy-101633464079686.html>.
- Pascanu R, Mikolov T, Bengio Y (2013) On the difficulty of training recurrent neural networks. Dasgupta S McAllester D, eds. *Internat. Conf. Machine Learn.*, Proceedings of Machine Learning Research, vol. 28(3) (PMLR, New York), 1310–1318.
- Paul RK, Saxena R, Chaurasia S, Zeeshan, Rana S (2015) Examining export volatility, structural breaks in price volatility and linkages between domestic and export prices of onion in India. *Agricultural Econom. Res. Rev.* 28:101–116.
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, et al. (2011) Scikit-learn: Machine learning in Python. *J. Machine Learn. Res.* 12(85):2825–2830.
- Pindyck RS (2004) Volatility and commodity price dynamics. *J. Futures Markets Futures Options Other Derivative Products* 24(11): 1029–1047.
- Pons N (2011) *Food and Prices in India: Impact of Rising Food Prices on Welfare* (Centre De Sciences Humaines, Delhi, India).
- Raka S, Ramesh C (2017) Understanding the recurring onion price shocks: Revelations from production-trade-price linkages. Policy paper, National Centre for Agricultural Economics and Policy Research, New Delhi, India, xvii + 56 pp.
- Raleigh C, Choi HJ, Kniveton D (2015) The devil is in the details: An investigation of the relationships between conflict, food price and climate across Africa. *Global Environ. Change* 32:187–199.
- Saha N, Kar A, Kumar P (2020) An investigation on performance of onion price and market arrivals in major Indian markets. *J. Agriculture Ecology* 9(9):78–82.
- Saxena R, Paul RK, Kumar R (2020) Transmission of price shocks and volatility spillovers across major onion markets in India. *Agricultural Econom. Res. Rev.* 33(1):45–52.
- Schumaker RP, Chen H (2009) Textual analysis of stock market prediction using breaking financial news: The AZFin text system. *ACM Trans. Inform. Systems* 27(2):1–19.
- Secomandi N (2015) Merchant commodity storage practice revisited. *Oper. Res.* 63(5):1131–1143.
- Shahaf D, Guestrin C (2010) Connecting the dots between news articles. *Proc. 16th ACM SIGKDD Internat. Conf. Knowledge Discovery Data Mining* (ACM, New York), 623–632.
- Sharma P, Gummagolmath K, Sharma R (2011) Prices of onions: An analysis. *Econom. Political Weekly* 46(2):22–25.
- Sugar CA, Gareth JM (2003) Finding the number of clusters in a data set: An information theoretic approach. *J. Amer. Statist. Assoc.* 98(463):750–763.
- Syntetos AA, Babai Z, Boylan JE, Kolassa S, Nikolopoulos K (2016) Supply chain forecasting: Theory, practice, their gap and the future. *Eur. J. Oper. Res.* 252(1):1–26.
- Tetlock PC (2007) Giving content to investor sentiment: The role of media in the stock market. *J. Finance* 62(3):1139–1168.
- The Hindu* (2015) Rice mill ‘scam’ triggers cag audit. <https://www.thehindu.com/news/national/Rice-mill-%E2%80%9898scam%E2%80%99-triggers-CAG-audit/article60330346.ece>.
- The Hindu Businessline* (2020) Onion price falls due to govt action against hoarding. <https://www.thehindubusinessline.com/economy/agri-business/onion-price-falls-due-to-govtaction-against-hoarding/article32940468.ece>.
- TOI (2008) Infected crop shoots up onion prices. <https://timesofindia.indiatimes.com/city/pune/Infectedcrop-shoots-up-onion-prices/articleshow/3910670.cms>.
- TOI (2009a) Transport strike, oil crisis push up prices. <https://timesofindia.indiatimes.com/city/kolkata/Transport-strike-oil-crisis-push-up-prices/articleshow/3958998.cms>.
- TOI (2009b) Onion floods Bhavnagar market. <https://timesofindia.indiatimes.com/city/rajkot/Onionfloods-Bhavnagar-market/articleshow/4030814.cms>.
- TOI (2009c) Andhra floods and strike fuel vegetable price rise in state. <https://timesofindia.indiatimes.com/city/bhubaneswar/Andhra-floods-and-strike-fuel-vegetable-price-rise-in-state/articleshow/5091545.cms>.
- TOI (2009d) Veggie prices spike on watery wave. (October 6), <https://timesofindia.indiatimes.com/city/bengaluru/Veggie-prices-spike-on-watery-wave/articleshow/5091855.cms>.
- TOI (2013a) Rice price plunges by 24. (August 12), <https://timesofindia.indiatimes.com/city/bengaluru/rice-priceplunges-by-24-on-bountiful-rain/articleshow/21770742.cms>.
- TOI (2013b) Rice export scam: Govt blocks CBI probe, 20 babus go scot-free. (December 24), <https://timesofindia.indiatimes.com/india/rice-export-scam-govt-blocks-cbi-probe-20-babusgo-scot-free/articleshow/27811928.cms>.
- TOI (2013c) Wheat exports from India to plummet as farmers hoard harvest. (May 22), <https://www.livemint.com/Politics/5Ect1st9e2OqQVr4pvAArO/Wheat-exports-from-India-to-plummet-as-farmers-hoard-harvest.html>.
- TOI (2013d) Railways to increase passenger fare, freight tariff by 2 per cent. (October 4), <https://timesofindia.indiatimes.com/business/india-business/railways-to-increase-passengerfare-freight-tariff-by-2-per-cent/articleshow/23533196.cms>.
- TOI (2013e) Protest against fuel price hike. (September 16), <https://timesofindia.indiatimes.com/city/varanasi/protestagainst-fuel-price-hike/articleshow/22619843.cms>.
- TOI (2014a) Monsoon delay fails to affect rice procurement. (August 5), <https://timesofindia.indiatimes.com/india/monsoon-delay-fails-to-affect-rice-procurement/articleshow/39651919.cms>.
- TOI (2014b) Paddy procurement scam probe yet to reach logical conclusion. (February 4), <https://timesofindia.indiatimes.com/city/patna/paddy-procurement-scam-probe-yet-to-reachlogical-conclusion/articleshow/29833580.cms>.
- TOI (2014c) HC asks state to explain rice ‘scam’. (February 12), <https://timesofindia.indiatimes.com/city/patna/hc-asksstate-to-explain-rice-scam/articleshow/30238593.cms>.
- TOI (2014d) Mamata Banerjee urges to check potato export. (August 8), <https://timesofindia.indiatimes.com/city/kolkata/mamata-banerjee-urges-to-check-potato-export/articleshow/39897577.cms>.
- TOI (2014e) Wholesale onion price dips to 5 per kg in Mumbai. (February 9), <https://timesofindia.indiatimes.com/business/india-business/wholesale-onion-price-dips-to-5-per-kg-in-mumbai/articleshow/30082568.cms>.
- TOI (2015a) FCI stuck with 24m tonnes of poor wheat. (June 29), <https://timesofindia.indiatimes.com/india/fcistuck-with-24m-tonnes-of-poor-wheat/articleshow/47857533.cms>.
- TOI (2015b) Poor monsoon results in lower rabi cultivation. (December 24), <https://timesofindia.indiatimes.com/city/nagpur/Poor-monsoon-results-in-lower-rabi-cultivation/articleshow/50303803.cms>.
- TOI (2015c) Government to give subsidy to traders to export potatoes. (March 6), <https://timesofindia.indiatimes.com/city/kolkata/government-to-give-subsidy-to-traders-to-export-potatoes/articleshow/46476011.cms>.
- TOI (2015d) 5 godowns raided, foodgrains worth over Rs 120cr seized. <https://timesofindia.indiatimes.com/city/thane/5-godowns-raided-foodgrains-worth-over-Rs-120cr-seized/articleshow/49500914.cms>.
- TOI (2015e) Raids continue, 930 quintals illegal dal stock seized. (October 23), <https://timesofindia.indiatimes.com/city/bhubaneswar/Raids-continue-930-quintals-illegal-dal-stock-seized/articleshow/49562079.cms>.

- TOI (2015f) Raids on dal mills, godowns continue. (October 26), <https://timesofindia.indiatimes.com/city/bhubaneswar/Raids-on-dal-mills-godowns-continue/articleshow/49533170.cms>.
- TOI (2015g) Truckers strike soar vegetable prices in Allahabad. <https://timesofindia.indiatimes.com/city/allahabad/truckers-strike-soar-vegetable-prices-in-allahabad/articleshow/49237710.cms>.
- TOI (2018) Day 4 of transporters' strike: Potato prices climb as supply falls by 40%. (October 6), <http://timesofindia.indiatimes.com/articleshow/65111423.cms>.
- Vaca CK, Mantrach A, Jaimes A, Saerens M (2014) A time-based collective factorization for topic discovery and monitoring in news. *Proc. 23rd Internat. Conf. World Wide Web* (ACM, New York), 527–538.
- Van Dijk TA (2013) *News as Discourse* (Routledge, New York).
- Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł, Polosukhin I (2017) Attention is all you need. Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, Garnett R, eds. *Adv. Neural Inform. Processing Systems 30 (NIPS 2017)* (Curran Associates, Inc., Red Hook, NY).
- Virk A (2015) History shows why no politician wants to mess with the onion. *The Quint* (August 20), <https://www.thequint.com/news/india/history-shows-why-no-politician-wants-to-mess-with-the-onion>.
- Wang Y, Agichtein E, Benzi M (2012) TM-LDA: Efficient online modeling of latent topic transitions in social media. *Proc. 18th ACM SIGKDD Internat. Conf. Knowledge Discovery Data Mining* (ACM, New York), 123–131.
- Weber R (2015) Welfare impacts of rising food prices: Evidence from India. Technical report, International Association of Agricultural Economists, Toronto.
- WorldGrain (2015) Weather issues drop India's 2014-15 grain production. (September 9), <https://www.world-grain.com/articles/5900-weather-issues-drop-india-s-2014-15-grain-production>.
- Wu D (2023) Text-based measure of supply chain risk exposure. *Management Sci.*, ePub ahead of print September 19, <https://doi.org/10.1287/mnsc.2023.4927>.
- Xiao G, Yang N, Zhang R (2015) Dynamic pricing and inventory management under fluctuating procurement costs. *Manufacturing Service Oper. Management* 17(3):321–334.
- Xie B, Passonneau R, Wu L, Creamer GG (2013) Semantic frames to predict stock price movement. *Proc. 51st Annual Meeting Assoc. Comput. Linguistics*, Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics, vol. 1 (Association for Computational Linguistics, Kerrville, TX), 873–883.
- Xing FZ, Cambria E, Welsch RE (2018) Natural language based financial forecasting: A survey. *Artificial Intelligence Rev.* 50(1):49–73.
- Yarlott WV, Cornelio C, Gao T, Finlayson M (2018) Identifying the discourse function of news article paragraphs. *Proc. Workshop Events Stories News 2018* (Association for Computational Linguistics, Kerrville, TX), 25–33.
- Zhang W, Skiena S (2010) Trading strategies to exploit blog and news sentiment. *Proc. Internat. AAAI Conf. Web Social Media* (AAAI Press, Palo Alto, CA), 375–378.
- Zhu X, Ninh A, Zhao H, Liu Z (2021) Demand forecasting with supply chain information and machine learning: Evidence in the pharmaceutical industry. *Production Oper. Management* 30(9):3231–3252.