

The Price is Right?

Statistical evaluation of a crowd-sourced market information system in Liberia

Anonymous
Anonymous
Anonymous
Anonymous
Anonymous

Anonymous
Anonymous
Anonymous
Anonymous
Anonymous

ABSTRACT

Many critical policy decisions depend upon reliable and up-to-date information on market prices. Such data are used to construct consumer price indices, measure inflation, detect food insecurity, and otherwise influence macroeconomic policy. In developing countries, where many of these problems are most acute, reliable market price information can be hard to come by. Here, we describe and evaluate Premise Data, a new technology for measuring price information using crowd-sourced data contributed by local citizens. Our evaluation focuses on Liberia, a fragile economy with a history of price insecurity, where Premise Data recently began collecting data. Our analysis utilizes tens of thousands of individual price observations collected at hundreds of different locations in Monrovia. We illustrate how these data can be used to construct composite market price indices, and compare these constructed indices to several sources of “ground truth” data from the Liberian Central Bank and the United Nations World Food Programme. Our results indicate that the crowd-sourced price data is strongly correlated with traditional price indices, but that statistically and economically significant deviations exist that require deeper investigation. We conclude by discussing how indices based on Premise data can be further improved with simple supervised learning methods that use traditional low-frequency data to calibrate and cross-validate the high-frequency Premise-based indices.

Categories and Subject Descriptors

J.4 [Computer Applications]: Social and Behavioral Sciences;
G.3 [Mathematics of Computing]: Probability and Statistics

General Terms

Human Factors, Measurement, Economics, Algorithms

Keywords

Consumer price index; premise; crowd-sourcing; ICTD

1. INTRODUCTION AND MOTIVATION

The Consumer Price Index (CPI) is one of the most important economic statistics used by policymakers to determine macroeconomic

policy and to evaluate the health of an economy [12]. It is the primary barometer for measuring inflation, which in turn impacts both fiscal and monetary policy); it is used to calculate purchasing power and determine exchange rates; it also directly impacts wages and welfare payment in both the private and public sector.

The CPI is intended to capture the overall cost of goods and services paid for by a typical individual at local markets. It requires two primary inputs: the “basket of goods” that is intended to be representative of a typical consumer; and the price for each of the goods in the basket. It is these prices that we focus on in this paper. In developed economies, consumer price data are typically obtained from a variety of sources, including supermarkets, service locations, and online retailers [4]. In the United States, for instance, “Bureau of Labor Statistics data collectors visit or call thousands of retail stores, service establishments, rental units, and doctors’ offices, all over the United States, to obtain information on the prices of the thousands of items used to track and measure price changes in the CPI.”¹ In developing economies, where most transactions are analog and the national statistical offices are more resource-constrained, there are fewer sources of accurate and up-to-date market price information [10].

The lack of reliable and up-to-date price information is particularly problematic in fragile economies, where dependence on subsistence agriculture and the lack of insurance and other social safety nets can exacerbate the impact on prices of weather shocks, political instability, and food insecurity. In Liberia and neighboring countries, for instance, food prices have been one of the primary instruments used in assessing the economic impacts of the recent Ebola outbreak [15, 6].

Here, we investigate a new technology for collecting price information in developing countries, which relies on “crowd-sourced” observations collected by local citizens with mobile phones. We focus on Premise Data, a technology platform that allows for mobile-equipped citizens to capture and upload price information to a central service. This technology, described in greater detail in Section 3.1, is similar to a small number of related platforms that enable crowd-sources data collection in developing countries.² The mClerk [7] and txtEagle [5] platforms are the most directly comparable systems of which we are aware; both use mobile-based platforms to gather data from low-end mobile phones. More broadly,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

¹http://stats.bls.gov/cpi/cpifaq.htm#Question_8, accessed July 2015

²A much more extensive literature, which we do not review here, explores the potential for data collection using mobile phones in developed economies. See [11] for an overview.

several examples exist of researchers sourcing data from the crowd: for instance, [14] study microblogging in response to a large earthquake in Haiti in 2010, and [1] describe several other emerging technological systems that facilitate participatory contribution in development areas including agriculture, rural development and natural resource management. In the closest study to our own, [8] conduct a feasibility study to explore the use of the Jana platform for collecting price data. While [8] demonstrate the potential of the platform, they focus on a description of the technological platform, rather than on evaluating the accuracy of the collected data in comparison to external sources of validation data.

Different from prior work, our focus is not on the technological artifact or interface used to collect data; rather, we study the data generated by this platform, and statistically evaluate its potential for use as an index of inflation and related economic activity. The empirical analysis relies on data generated from the Premise network of contributors in Liberia, which has been collecting price information on 38 market goods in Monrovia. We present the dataset and document several prominent features related to the stability and noise present in the raw data, and discuss potential corrections and processes for removing outliers that may improve the reliability of derivative metrics (Section 3.2). We also describe a basic method for computing a consumer price index from the crowd-sourced data (Section 3.3), and compare these indices to alternative measures of inflation in Liberia. Finally, we discuss preliminary efforts to improve the accuracy of price indices produced from crowd-sourced Premise data by combining the cleaned Premise data with data from traditional, low-frequency sources.

This paper thus makes three primary contributions. First, we illustrate how crowd-sourced market data can be used to construct price indices, and characterize the statistical properties of these indices. Second, we carefully evaluate these indices by comparing them to several more traditional methods for measuring food insecurity. Finally, we provide a simple supervised learning framework that can be used to improve the accuracy of high-resolution estimates through calibration and cross-validation with low-resolution sources of data.

2. BACKGROUND: CONSUMER PRICE INDICES IN LIBERIA

Following Liberia’s long-entrenched civil war and post-conflict recovery, which ended in 2003, the country took several years to achieve political and economic stability. Between 2011 and 2014, consumer prices in Liberia followed a fairly linear trend with mean year-to-year inflation at approximately 9%. In the period following the Ebola Virus Disease epidemic, the country saw an uptick in the price index, with year-to-year inflation rising to a peak of 16.3% in September 2014. With the rise in consumer prices, concern over access to food was heightened. As of April 2015, the government and international agencies remained vigilant to the threat to food security that increased prices posed for Liberian households [?].

Price data is historically accessible to a limited set of actors and on a limited basis. Official price data is collected by the Liberian government and analyzed by the Central Bank of Liberia in order to produce aggregate price indices on a monthly basis. More frequent data collection (collected by the World Food Programme) is typically collected for the purpose of monitoring food insecurity. However, data collected for food security monitoring, understandably, covers a more limited set of consumption items and locations relative to the government statistics.

Price data plays an integral roll in the WFP’s forecasts of food shortages, as Food prices influence household spending decisions. As such, the WFP uses food prices as an indicator of the impact of economic shocks on households. Access to timely price data has been especially important during the 2014-15 Ebola epidemic in West Africa, where the combination of restricted economic activity, loss of life, and loss of sources of income have culminated in a threat to food security [?].

3. PREMISE DATA

3.1 Technology platform

Premise (www.premise.org) is a technology company based in San Francisco that is developing a platform for capturing data from a distributed network of individual contributors. The intent of the platform is to enable rapid and adaptive measurement of local economic and social infrastructure, using data collected by local citizens. Premise recruits individuals in urban and rural regions of developing countries to perform simple, structured tasks that capture information about their local community (Figure 1). Premise currently operates in 30 countries worldwide. A major focus of Premise’s efforts to date has been on collecting price data from developing countries.

Premise contributors use photo-enabled phones to capture prices in local markets. Contributors are compensated in the local currency, and are trained in person prior to submitting data. Each day, contributors receive a list of *tasks* which detail the items for which price observations are needed. Contributors are typically paid piecemeal for each successfully completed task, for an amount “on the order of the price of an egg” for each data point captured.³ Photos and price information submitted by the contributor go through a quality control screening process, described in more detail below.

Beginning in November 2014, Premise initiated data collection in Liberia in an effort to produce a Food Staples Price Index (FSPI) to collect and other measures of economic and business activity. The initial focus has been in Monrovia, where Premise is currently collecting daily price observations for staple foods and non-food items in all of Monrovia’s major market areas (Figure 1). In the summer of 2015, Premise plans to expand data collection across the country, beginning with Voinjama and Fish town in July 2015.

Nine months after initiating data collection in Monrovia, Premise receives approximately 600 price observations per week for a basket of 38 unique products, from a small network of independent contributors. As can be observed in Table 1, which presents summary statistics for the raw contributor data, there is a great deal of variation in the frequency at which each product is observed, the number of unique locations at which a product is captured, and the price level and variance for each product over time. As we describe in the following section, the number of errant observations also varies considerably by product.

3.2 Detecting outliers in crowd-sourced data

The raw data captured by Premise contributors in Liberia are illustrated in Figure 2. Here, we plot a separate time series for each of the 38 products as a semi-transparent grey line. Each of these lines represents the daily average price for that product, averaged across all observations taken on that day by all contributors in all locations. Highlighted in blue are the 6 time series corresponding to products in the “Grains and flours” category. In red is the

³Based on private correspondence with Premise staff, July 2015.

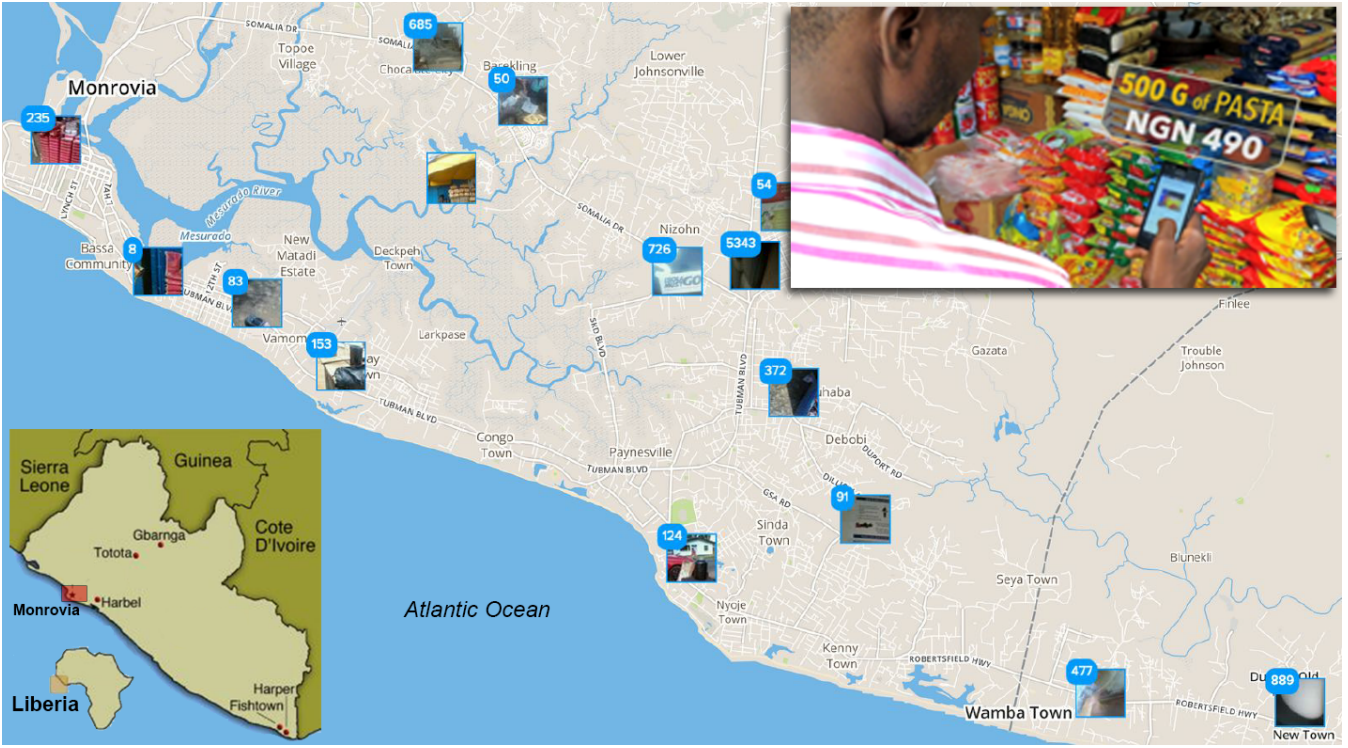


Figure 1: **Premise Data collection methodology.** Premise indexes and analyzes data captured by a global network of contributors. Bottom-left: Location of Liberia and its capital Monrovia. Main figure: Locations from which contributors have captured data. Numbers indicate the number of data points collected from each location over the past three months; square icons are actual images uploaded by contributors. Top-right: Schematic of data capture process in which a contributor uses a cameraphone to photograph the prices of pasta at a local market. These photos are sent to Premise and form the basis for the data we analyze.

composite sub-index, calculated using a procedure we will shortly describe.

As is evident in the product-level time series in Figure 2, the raw data collected by Premise contributors is subject to several sources of error. On some occasions there appear to be idiosyncratic spikes where a single product’s price will change by as much as 200%; at others, these spikes appear to be correlated across products. Of primary concern is disentangling from actual changes in prices from measurement error. As has been documented in related work, there are many possible sources of such error, both accidental and deliberate [8, 3, 2]. These include input errors (for instance, a misplaced decimal point or a photograph of the floor), as well as outright fraud where a contributor intentionally falsifies data. Premise implements several measures to detect and prevent such deliberate fraud [13], but many of these are not publicly disclosed, and in practice affect a relatively small number of total captured data.

In follow-up work, we are developing more refined techniques for identifying and removing erroneous data, which may constitute as much as 20% of the total data captured on the Premise platform. Here, we describe a simple procedure that, based on manual verification, appears to catch a large share of these errors. Formally, we denote by P_{itlk} an observation recorded by individual i for item k in location l at time t . We define price outliers as those observations that deviate significantly from historical prices for a given product, i.e.,

$$|\log(P_{itlk}) - \log(\mu_{ilk})| > \lambda_k \sigma_{ilk} \quad (1)$$

where $\mu_{ilk} = \frac{1}{Nn} \sum_i \sum_{s < t} P_{islk}$ is the average historical value for k at location l (assuming N individuals and n observations where $s < t$), and σ_{ilk} is the corresponding standard deviation. In this framework, λ_k is the key parameter which determines the stringency with which outliers will be identified. In practice, Premise currently applies a common threshold across all products of $\lambda_k = \lambda \approx 2$. In ongoing work, we are exploring a supervised learning approach to determining a product-specific λ_k , which will allow for some products with greater expected variation over time to exhibit more intertemporal variability.⁴

3.3 Computing CPI from crowd-sourced data

A primary objective of the Premise application is to convert the disaggregated price data collected by the network of contributors into more meaningful price indices, similar to the CPI, which can then be used to measure inflation and inform food policy decisionmaking. Here, we describe the process used to construct the Food Staples Price Index (FSPI), the Premise equivalent of the CPI, from the disaggregated data. Formally, our goal is to compute an aggregate CPI_{rm} for a region r in month m , from a large number of disaggregated price observations P_{itlk} .

Given the set of P_{itlk} with outliers removed, the FSPI CPI_{rm} is constructed using a process based loosely on the methods employed by the US Bureau of Labor Statistics [4]. Initially, the average daily

⁴Similarly, when insufficient observations exist from which to derive reliable estimates of μ_{ilk} and σ_{ilk} , a “bootstrap” process is used to manually curate and reject anomalous observations.

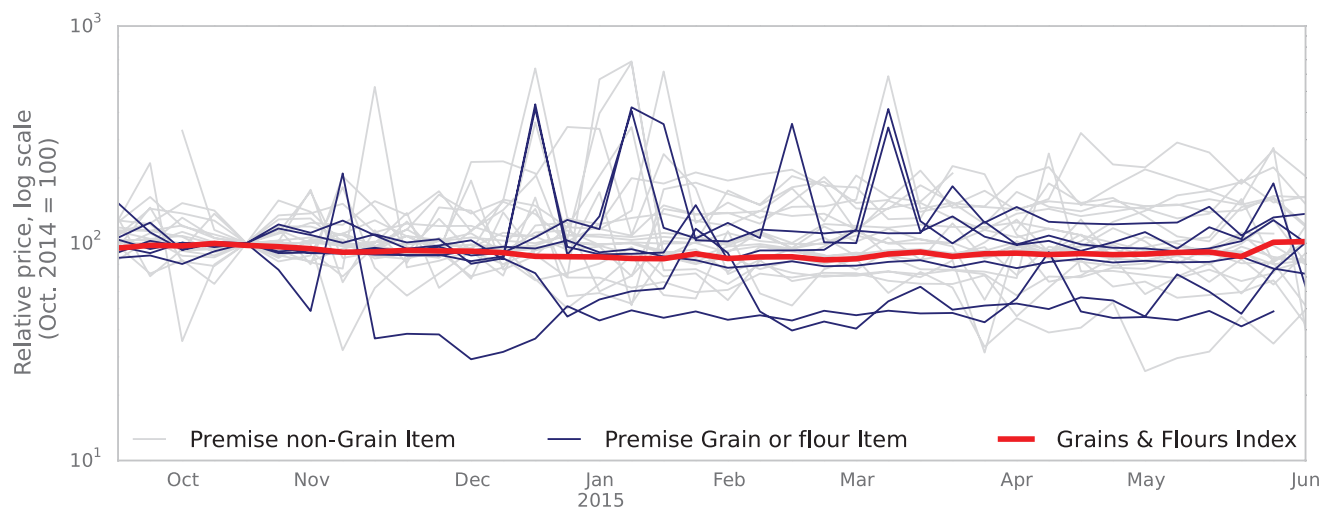


Figure 2: **Raw data received from Premise contributors.** Data for 38 different products is captured by Premise contributors (light grey lines). Of these goods, 6 are in the “Grains and Flours” product group (dark blue lines), corresponding to bread, bulgur wheat, butter rice, cassava flour, fan-fan rice, and USA parboiled rice. Using the methods described in Section 3.3, these six product price-series are aggregated into a single sub-index for the product group (red line). Each series is shown normalized to a value of 100 on October 15, 2015.

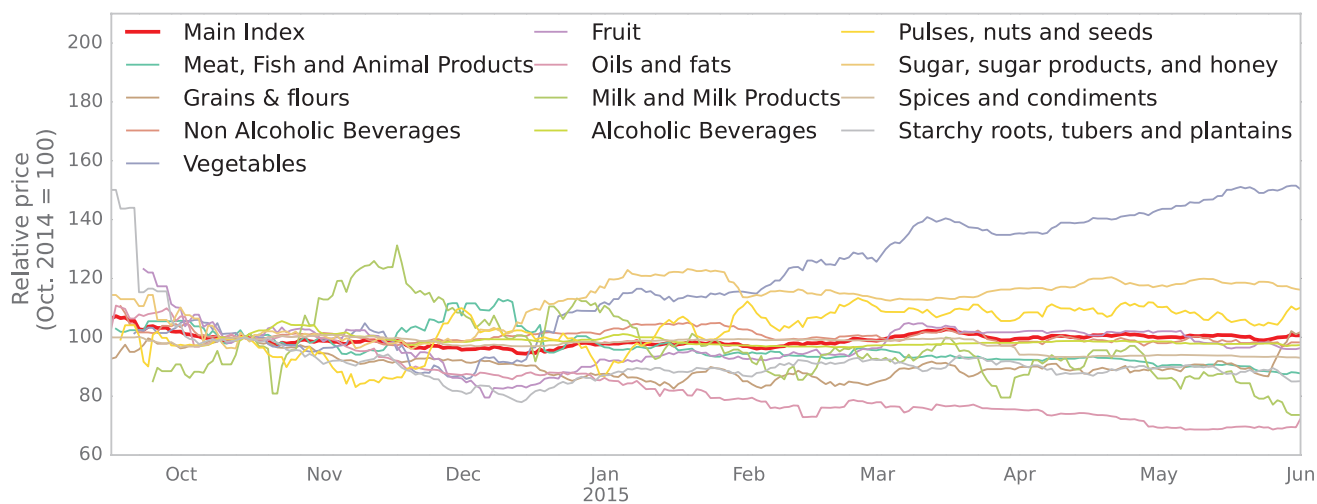


Figure 3: **Premise FSPI and sub-indices, as calculated from contributor data.** Following the procedure described in Section 3.3, sub-indices are computed for each of the 12 product groups listed in Table 2. Using the weights listed in the same table, the composite Food Staples Price Index, the Premise equivalent of a CPI, is calculated and shown as a thick red line.

Index component	weight (%)
Alcoholic Beverages	1.2%
Fruit	9.2%
Grains and Flours	13.7%
Meat, Fish and Animal Products	16.8%
Milk and Milk Products	6.1%
Non Alcoholic Beverages	12.0%
Oils and Fats	3.0%
Pulses, Nuts and Seeds	2.9%
Spices and Condiments	8.0%
Starchy Roots and Tubers	11.4%
Sugar and Sugar Products	2.2%
Vegetables	13.4%

Table 2: **Weights used in constructing the FSPI.** The composite Food Staple Price Index is constructed as the weighted sum of 12 primary sub-indices, where the above weights are determined based on a recent household expenditure survey.

price P_{Tlk} for item k at location l is constructed by taking the average of all $|T|$ observations collected on day T , i.e.,

$$P_{Tlk} = \frac{1}{N|T|} \sum_{i \in l} \sum_{t \in T} P_{itlk} \quad (2)$$

This value is still quite specific, as l can be as precise as a single storefront location, and k can be unique to an item SKU, such as the price of one bottle of Club Beer (a local beer), or the cost of a bucket of low-grade gari (a local flour). These item-day averages are next aggregated into product-day averages P_{TIK} by standardizing units of measurement (e.g., pounds to grams) and taking the geometric mean of related products (e.g., beef briskets), i.e.,

$$P_{TIK} = \left(\prod_{k=1}^K P_{Tlk} \right)^{1/K} \quad (3)$$

Product-day averages are similarly aggregated across all locations in a region and across all days in a time window to produce monthly estimates of the regional cost of a specific product. These product averages are further aggregated into product sub-groups (e.g., standard cut beef), product groups (e.g., beef), and sub-indices (e.g., meat). This aggregation uses weights that are determined based on the estimated expenditures of consumers on the various items. In Liberia, these weights are determined by a recent consumer expenditure survey [?].

In Liberia there are 12 such *sub-indices*, which indicate the prices of the most common items in the country. The final step in constructing the FSPI is to combine the sub-indices into a single consumer price index that reflects the price level of a typical market basket of food staples. The weights w_k for each of the sub-indices used in constructing the Liberian FSPI are given in Table 2. Thus, the composite FSPI can be expressed as a weighted aggregate of the original contributor observations:

$$CPI_{rm} = \frac{1}{|r||m|} \sum_{l \in r} \sum_{t \in m} w_k P_{Tlk} \quad (4)$$

4. EVALUATION RESULTS

The methods described above make it possible to construct a CPI-like metric, the FSPI, as well as several sub-indices of product-group prices, from the data captured by Premise contributors. In order to validate the relevance of these constructs, we take two approaches - one at the aggregate index level and another at the item level. We draw from two sources of data for the validation. The purpose of this comparison is to judge the consistency and reliability of Premise data vis-a-vis an economic indicator for the Liberian economy as well as a best-available option for policymakers interested in food security as proxied by the price of individual goods. The data for the former comparison comes from the Central Bank of Liberia, while the item-level comparison is conducted with a primary source of data collected for the United Nations World Food Programme with the purpose of tracking threats to food security in Liberia.

4.1 Comparison data

4.1.1 Central Bank of Liberia

The Central Bank of Liberia releases headline consumer price index data from the 15th day of each month. National price level are based on monthly price surveys conducted by the Liberia Institute of Statistics and Geo-information Services (LISGIS). Price indices are provided for the overall price level, food and non-alcoholic beverages (split by domestic and imported food), transportation, and imported fuel. The inset of Figure 4 displays the time series of the Food and Non-Alcoholic Beverages Index for Liberia between May 2011 and April 2015. To emphasize relative changes in prices, the index is normalized so that the value of the index in October 2014 is equal 100. The impact of the Ebola epidemic, which caused year-to-year price inflation to peak at 16% in September 2014, is evident in the figure. Note that the Central Bank index is intended to be nationally representative, but Premise data for this period is restricted to the capital city, Monrovia.⁵

4.1.2 World Food Programme

We additionally compare the product-level data captured by Premise contributors to data acquired from United Nations World Food Programme (WFP), which conducts regular data collection for key food prices in order to assess food security and identify price shocks that disproportionately affect poor households. We utilize data collected for the WFP on a key set of consumption goods in Liberia, which were initiated in part to monitor the impact of Ebola on price inflation. For the purpose of comparability with the Premise data, we restrict the WFP data to a subset of four goods: imported rice, cassava, palm oil and charcoal.⁶ Figure 5 shows the time series of mean weekly price of both rice and cassava, as independently observed in Premise and WFP data.

4.2 Statistical comparison

Comparison to baselines

We compare the Premise data to both the Central Bank food and non-alcoholic beverages index and the WFP individual item prices for the period from September 2014 to April 2015. Effectively, we seek to quantify the differences in Figures 4 and 5. These results are presented in Table 3, which indicates the average per-month error (RMSE) as well as the correlation between datasets

⁵In July 2015, Premise initiated data collection in two rural markets, but these data were not available at the time this research conducted.

⁶The WFP data captured six products in total, including brown cowpeas and gari flour; however, the Premise data collection does not include these items.

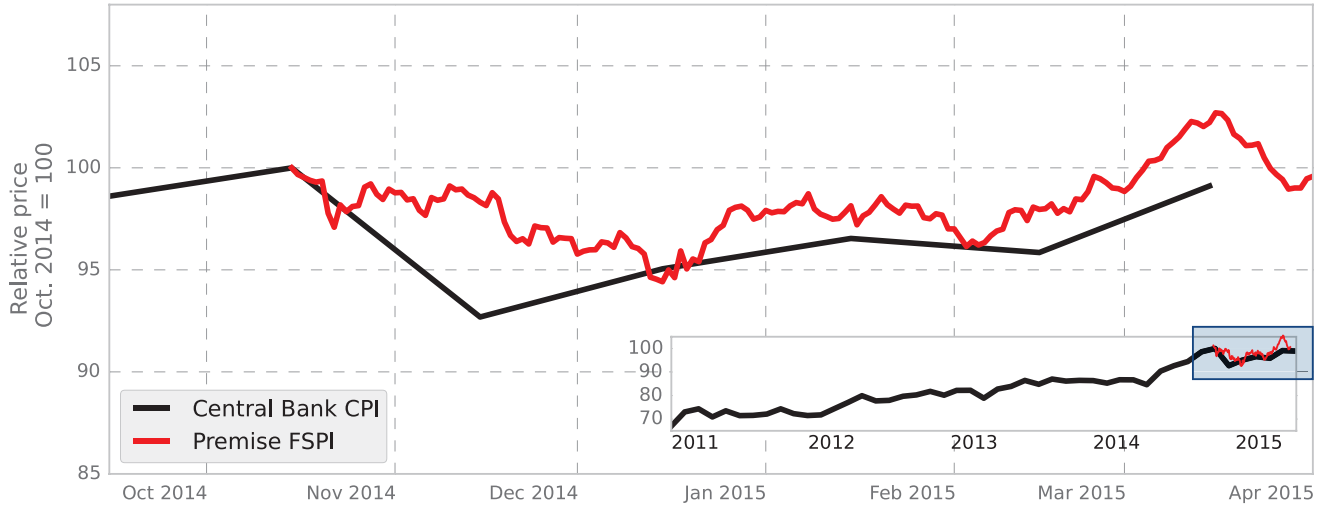


Figure 4: **Comparison of FSPI to Central Bank CPI.** Main figure shows the CPI calculated by the Central Bank of Liberia and the FSPI based on Premise data. Inset figure displays the Food and Non-Alcoholic Beverages relative price from May 2011 to April 2015.

over time. Since the CPI data is collected at the monthly level while the Premise data exists in daily averages, we compare the two series by making monthly comparison of the CPI to the average Premise data over the preceding 30-day period (row 1), and also by linearly interpolating the CPI data between monthly observations and comparing at the daily level (row 2).

As is evident in the Table, there is a relatively strong correlation between the aggregate indices, but the correlation is weaker when analyzing specific products. Importantly, these correlations do not account for potential delays and offsets in the different series. For instance, if the Premise FSPI is a leading indicator of the CPI, or vice versa, such patterns would not be reflected in the results in Table 3.

Modeling improvements

Our analysis thus far indicates a relatively strong correlation between the premise data and the indices collected by the Liberian Central Bank and the UN WFP. Moving forward, we believe a promising area for research lies in understanding how the multiple data sources can be reconciled into meta-indices that reflect variation in prices at different levels of temporal resolution. While a comprehensive study of these possibilities is beyond the scope of the current work, we briefly describe one such line of research that we believe would be fertile ground for future work.

Specifically, we are interested in understanding the extent to which future predictions of inflation, currently based on Central Bank CPI data, might be improved with the disaggregated and high-frequency data collected by Premise. Thus, we begin with a simple time-series (ARIMA) model, which provides an indication for how well the CPI can be forecast, relying only on historical CPI data. We compare this to two (vector autoregression) alternatives, that additionally incorporate historical data from (i) the Premise FSPI, and (ii) all of the Premise sub-indices shown in figure 3.

To construct a baseline, we use the Hyndman-Khandakar algorithm for automatic ARIMA model selection. This algorithm uses repeated KPSS tests to determine the number of times the data must

be differenced to be stationary.⁷ A stepwise search over values of p and q minimizes the AICc. For the full four-year panel of CPI data, the best model is ARIMA(0,1,1) with drift (AICc=215.4); over the 9-month period for which the CPI data intersects with the Premise data period, a “white noise” ARIMA(0,0,0) model with non-zero mean is selected (AICc=42.64).

We compare the baseline to a vector autoregression (VAR) model that provides structure to capture the relationships between CPI and the overlapping time-series data from Premise [9]. Formally, if $Y_t = (y_{1t}, \dots, y_{kt})$ denotes a vector of k time series, we use the VAR model

$$Y_t = c + A_1 Y_{t-1} + A_2 Y_{t-2} + \dots + A_p Y_{t-p} + \varepsilon_t \quad (5)$$

where A_i are $(k \times k)$ coefficient matrices and ε_t is a white noise vector process.

Estimation results are presented in Table 4. We observe that the combined VAR models, which take advantage of the variation in Premise data as well as the historical CPI data, perform marginally better at forecasting than models based solely on the historical CPI. While this framework is rather rudimentary, there is preliminary evidence that such an approach might provide an avenue toward developing a set of price indices using data from multiple sources.

5. DISCUSSION

While the Premise FSPI is thus correlated with the two traditional price measures we were able to obtain, there are statistical discrepancies that are not easily resolved.⁸ While we largely treat the Central Bank’s CPI data as “ground truth” and assume that deviations

⁷Separately, unit root tests indicate that the full panel of CPI data requires first-differencing, but the 9-month intersecting period does not need differencing; we therefore use first-differencing in the AR model and no differencing in the VAR models.

⁸One obvious source of these difference may be the differences in sampling frames used in data collection, for instance the fact that the Central Bank surveys the entire country’s prices while Premise is thus far focused in Monrovia; as Premise expands to additional markets it will be possible to test this hypothesis.

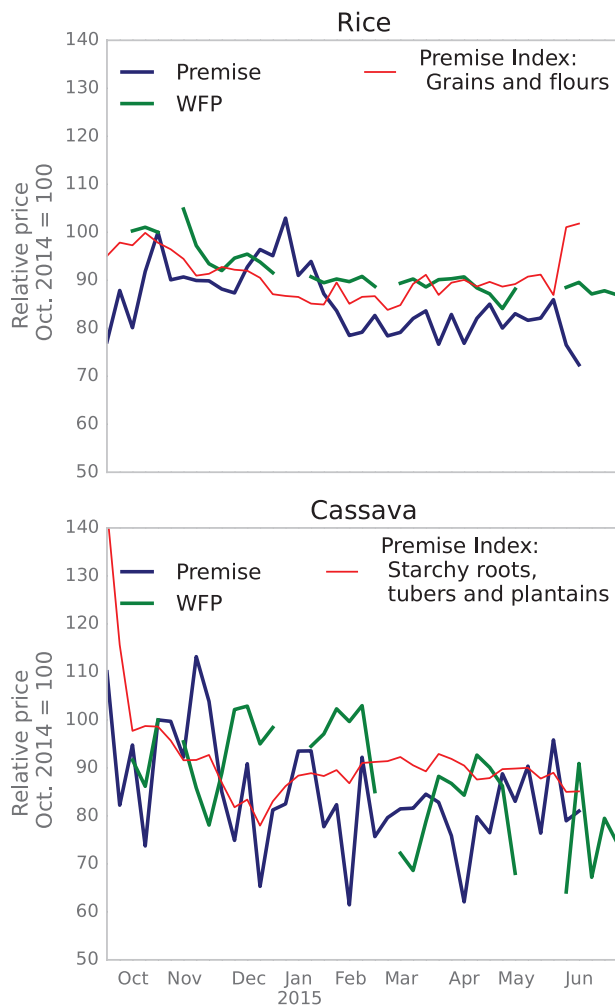


Figure 5: **Comparison of Premise and WFP Item Prices.** Price data for two food items in Liberia. The blue lines represent the average daily price data collected for rice (top figure) and cassava (bottom figure). The green lines indicate the corresponding price data for the same item, collected by the World Food Programme's Building Markets initiative. The red line indicates the composite sub-index, constructed from all goods in the grains (top figure) and starchy roots (bottom figure) product categories.

observed in the Premise FSPI are errors, it is also conceivable that the Premise data might at times be accurate where the Central Bank data is not. Indeed, the two sources of data, while comparable, have distinct advantages and disadvantages that we discuss briefly before concluding.

5.1 Advantages of traditional price data

A primary advantage of traditional sources of price information is that they are, in a word, traditional. Governments and international organizations have well-established mechanisms for collecting and processing price data, and a standard set of best practices exist for determining sample frames, deciding price frequencies, and integrating the resulting measures into macroeconomic decisionmaking. Centralized administration of these efforts further ensures the financial stability and continued support for data collec-

tion, whereas private sector efforts may be more subject to changing business models and sources of revenue. Of course, centralized government administration also makes it easier for published indices to be manipulated and coerced.

There are also strong economies of scale in centralized data collection. Governments and international organizations conduct data collection, such as censuses, expenditure surveys and firm surveys, that are complimentary to price data. For example, expenditure surveys are integral to updating consumption basket estimates that feed into consumer price indices.

5.2 Advantages of Premise CPI

Relative to traditional models of price data collection, the Premise platform offers several distinct advantages:

- **Granularity:** Premise data can be sourced continuously in time and space, increasing the ability to track prices in real-time, in sub-regions of the country.
- **Flexibility and Scalability:** While our focus has been on price data collection, the crowd-sourcing framework can be used to capture a much wider array of data types. Near-term possibilities include street mapping, collecting information on the availability of public utilities, and mapping financial inclusion. Because of the way that contributors are sourced and compensated, it is possible to quickly scale up data collection efforts.
- **Transparency:** Premise data is currently publicly distributed at www.premise.org, facilitating the use of the data in policymaking and research.

6. CONCLUSIONS

We describe and evaluate Premise data, a platform for collecting crowd-sourced price information from networks of local contributors in developing countries. Our focus on the statistical properties of the Premise data, and on comparing Premise-based indices to more traditional measure of price inflation, reveals several promising areas for future work. First, further quantitative work would benefit greatly from a longer panel of price data. While the Premise data contains tens of thousands of observations and is collected at extremely high frequency, the authoritative central bank data is collected only monthly; with only 9 months of overlapping data, it is very difficult to make robust statistical comparisons, and our initial attempts at time series modeling quickly become rather futile.

Second, there is considerable scope for improvement in the techniques used to identify erroneous data, caused both by innocent error and intentional fraud. The outlier removal system we describe and implement appears to be reasonably effective, but is rather coarse and relies heavily on (possibly unjustified) intuition. Given a labeled training set, where the source of erroneous data points is known, it would be possible to develop far more sophisticated methods that are potentially specific to a given product, location, or contributor.

Finally and perhaps most promising, we believe significant progress can still be made in developing methods for supervised learning that use authoritative data to improve the accuracy of estimates based on high-frequency data from Premise and related sources. Here, we made the simple point that inflation forecasts appear to improve when historical CPI data is supplemented with data from

Premise. Analogous approaches could be used to increase the granularity of official CPI estimates (beyond the country-month), or to construct monthly Premise estimates that correspond more closely to official benchmarks. As before, the absence of a long panel of training data makes these exercises difficult in the immediate term, but as data from Premise’s global network continues to stream in, it will open many opportunities for research that can impact how prices and inflation are measured in developing economies.

	R^2	RMSE
ARIMA (4 years)	0.23	2.47
ARIMA (9 months)	0.25	2.33
VAR (FSPI only)	0.15	3.09
VAR (FSPI+CPI)	0.24	2.01
VAR (FSPI+CPI+sub-indices)	0.31	1.89

Table 4: **Predictive power.** Comparison of forecasting performance of autoregressive (ARIMA) and vector autoregression (VAR) models. Whereas the ARIMA model incorporates only historical CPI data, the VAR model also utilizes data from Premise.

Market good	# observations	# locations	Units	Min	Max	Mean	SD
Banana	819	201	piece	0.33	25	10.47	3.73
Beef Brisket (fresh, raw)	317	102	piece	5	200	42.42	37.30
Beer	718	177	cl	0	1833.33	11.56	74.65
Bitter Ball	476	130	kg	0	475	74.28	53.98
Bitter Kola	823	192	piece	0.91	35	13.96	4.7
Boiled Eggs	891	206	piece	1	25	13.24	2.54
Bread	1002	202	piece	5	500	72.35	48.49
Bulgur Wheat	830	168	kg	0.03	3187.5	49.91	116.3
Butter Rice	886	183	kg	0.68	7437.5	85.2	353.04
Cassava	749	165	piece	0.11	100	15.12	9.12
Cassava Flour	979	203	kg	0	4250	57.45	228.58
Cassava Leaf	377	109	piece	8.33	100	18.90	14.07
Charcoal	807	193	kg	2.4	1190	17.25	59.93
Chloride	842	199	liter	0	2000	152	88.76
Fan-fan Rice	839	170	kg	1	4250	68.77	147.18
Fuel (Diesel)	569	164	liter	10.57	1733.33	100.95	111.06
Instant Coffee	518	154	gram	0.01	10	2.72	2.63
Kidney Beans (dry)	805	175	kg	0.18	15937.5	203.45	869.62
Live Chicken (medium size)	361	81	piece	0	1300	615.95	244.02
Mayonnaise	613	168	ml	0.01	739.34	35.02	96.95
Onion	1098	214	piece	0.4	125	21.47	20.59
Orange	852	193	piece	2.5	125	16.57	11.29
Palm Butter	711	158	kg	3	2500	35.37	149.35
Palm Oil	997	198	liter	9.51	1800	110.71	80.40
Petrol (gas)	682	175	gallon	0.31	710	252.43	111.52
Plantain (Cooking Banana)	900	198	piece	0	250	20.37	11.26
Potato Greens	608	141	piece	0.5	35	14.2	4.74
Powdered Milk	846	197	gram	0	70	0.86	2.78
Salt	918	190	gram	0	4.25	0.09	0.23
Sardines (canned)	332	99	gram	0.04	4.25	1.55	1.53
Seasoning cube	434	108	gram	0	127.5	0.72	6.11
Smoked Fish	592	162	piece	0.52	550	70.24	81.93
Sugar	967	198	gram	0	8.75	0.12	0.54
Tomato	564	132	piece	0.08	500	22.07	33.62
USA Parboiled Rice	981	205	kg	5	6250	92.21	279.91
Vegetable Oil	1012	197	liter	1.35	580	130.4	41.09
Water (Bag)	1035	217	ml	0	5.71	0.02	0.18

Table 1: **Summary statistics of Premise contributor data.** Data collection in Liberia began in October, 2014, with products and market locations being gradually added since then. There is significant variation between products in terms of average prices, price variation over time, presence of outliers, and unit measurements.

	Corr (A)	MSE (B)	MSE (% of mean) (C)
A. Central Bank Food Price Index			
FSPI (Monthly Average)	0.654	2.37	5.82%
FSPI (Daily vs. Central Bank linear interpolation)	0.654	2.10	4.56%
B. Food Items (WFP, relative prices)			
Rice	0.541	8.94	87.40%
Cassava	0.033	16.27	304.61%
Charcoal	0.351	6.29	41.70%

Table 3: **Model performance.** Measures of model accuracy and error, comparing Premise data to data collected by the Liberian Central Bank and the World Food Programme.

7. REFERENCES

- [1] H. Ashley, J. Corbett, D. Jones, B. Garside, and G. Rambaldi. Change at hand: Web 2.0 for development. *Participatory Learning and Action*, 59(1):8–20, 2009.
- [2] B. Birnbaum, G. Borriello, A. D. Flaxman, B. DeRenzi, and A. R. Karlin. Using Behavioral Data to Identify Interviewer Fabrication in Surveys. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '13, pages 2911–2920, New York, NY, USA, 2013. ACM.
- [3] B. Birnbaum, B. DeRenzi, A. D. Flaxman, and N. Lesh. Automated quality control for mobile data collection. In *Proceedings of the 2nd ACM Symposium on Computing for Development*, page 1. ACM, 2012.
- [4] BLS. The Consumer Price Index. In *Bureau of Labor Statistics Handbook of Methods*, number 17. Dec. 2008.
- [5] N. Eagle. txteagle: Mobile Crowdsourcing. In N. Aykin, editor, *Internationalization, Design and Global Development*, number 5623 in Lecture Notes in Computer Science, pages 447–456. Springer Berlin Heidelberg, 2009.
- [6] R. Glennerster and T. Suri. Economic Impacts of Ebola: Bulletin Four. *IGC Bulletin*, May 2015.
- [7] A. Gupta, W. Thies, E. Cutrell, and R. Balakrishnan. mClerk: Enabling Mobile Crowdsourcing in Developing Regions. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '12, pages 1843–1852, New York, NY, USA, 2012. ACM.
- [8] N. Hamadeh, M. Rissanen, and M. Yamanaka. Crowd-sourced price data collection through mobile phones. 2013.
- [9] J. D. Hamilton. *Time Series Analysis*. Princeton University Press, Princeton, N.J, 1 edition edition, Jan. 1994.
- [10] M. Jerven. *Poor numbers: how we are misled by African development statistics and what to do about it*. Cornell University Press, 2013.
- [11] N. Lane, E. Miluzzo, H. Lu, D. Peebles, T. Choudhury, and A. Campbell. A survey of mobile phone sensing. *IEEE Communications Magazine*, 48(9):140–150, Sept. 2010.
- [12] N. Mankiw. *Principles of Macroeconomics*. Cengage Learning, Jan. 2014.
- [13] Premise. Premise Data Corporation: Food Staples Indexes, Nov. 2014.
- [14] K. Starbird and L. Palen. "Voluntweeters": Self-organizing by Digital Volunteers in Times of Crisis. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, CHI '11, pages 1071–1080, New York, NY, USA, 2011. ACM.
- [15] World Food Programme. Liberia (April 2015): Rice and oil prices rise significantly in Lofa County. *mVAM Bulletin*, Apr. 2015.