



# **Big Data Analytics and Hadoop**

## **24ASD511**

*A case study on*  
**Food price analysis in India**

*Submitted By*

**Additi CB.PS.P2ASD24001**  
**Sivapriya CB.PS.P2ASD24025**

*Dataset*  
Food prices India - Kaggle

*Platform*  
Hive (on WSL)

# Index

S. NO	TITLE	PAGE NO
	<b>Abstract</b>	i
1	<b>Introduction</b>	1
2	<b>About the Tool</b>	2
3	<b>Objective</b>	2
4	<b>Architecture Diagram</b>	3
5	<b>Pathway</b>	4
6	<b>Dataset Description</b>	5 - 6
7	<b>Questions Overview</b>	7 - 8
8	<b>Data Loading</b>	9 - 10
9	<b>Data Analysis</b>	11 - 31
10	<b>Conclusion</b>	32 - 33
11	<b>References</b>	34

## **Abstract**

Big Data Analytics is transforming how organizations extract insights from the massive data generated daily. Probability and Statistics (ProbStat) are crucial in this domain, offering the tools needed to analyze patterns and make informed predictions, vital for applications like market analysis and financial forecasting. This case study explores Indian food prices using Big Data techniques to reveal trends and insights. By applying advanced SQL methods such as window functions and grouping sets, we analyze commodity price changes across different regions and seasons. ProbStat's role in real-time analytics allows businesses to address uncertainties, leading to proactive strategy adjustments. The key outcomes include actionable insights for pricing strategy optimization. The study utilizes Apache Hive for processing and SQL for query execution, with visualizations highlighting the results. This work demonstrates the importance of Big Data Analytics and ProbStat in turning data into strategic assets, fostering innovation and efficiency in a data-driven world.

## **Introduction**

This research focuses on the problem of typifying and forecasting the cash determinants within a specific segment of the economy – the Indian food marketplace. Price fluctuations, different commodities, and diverse geographical market areas of the country render these phenomena very complex and difficult to understand. An effective and complex supply chain management can be established only if adequate and relevant business intelligence is integrated. This study hopes to utilize the modern Big Data methodologies for in-depth analysis of the patterns of commodity prices and arm the relevant stakeholders with proper information.

The social relevance of this study is profound, as pricing affects the life of billions from farmers to consumers. Better understanding of price fluctuations can help stabilize the market and mitigate food waste while guaranteeing reasonable prices to be paid and profit earned. Reliable information on market pricing and tracking of volatile prices provided by this study helps narrow the resource gap and enhance the efficiency in the provision of basic food commodities. Improved endeavours will significantly increase food security and economic welfare, and therefore, this research is important for our society.

In order to gain these insights, the study uses tools like Apache Hive for data preprocessing and complex query execution and Matplotlib (Google CoLab) for visualizing the hive results. Commodity prices are analyzed using advanced SQL methods, like CTE's, window functions and grouping sets, to discover trends and correlations. The objective is to formulate insights that enable stakeholders to make decisions that improve pricing and market competitiveness. The study is designed in such a way that it helps companies to identify patterns and trends so they could better monitor and respond to changes, thus making the food market more stable and fair.

## About the Tool

For this study, Apache Hive has been selected as the primary tool for data processing and analysis. The Apache Hive is a distributed, fault-tolerant data warehouse system that enables analytics at a massive scale and facilitates reading, writing, and managing petabytes of data residing in distributed storage using SQL. Hive Metastore(HMS) provides a central repository of metadata that can easily be analyzed to make informed, data driven decisions, and therefore it is a critical component of many data lake architectures. Hive is built on top of Apache Hadoop and supports storage on S3, adls, gs etc though HDFS.

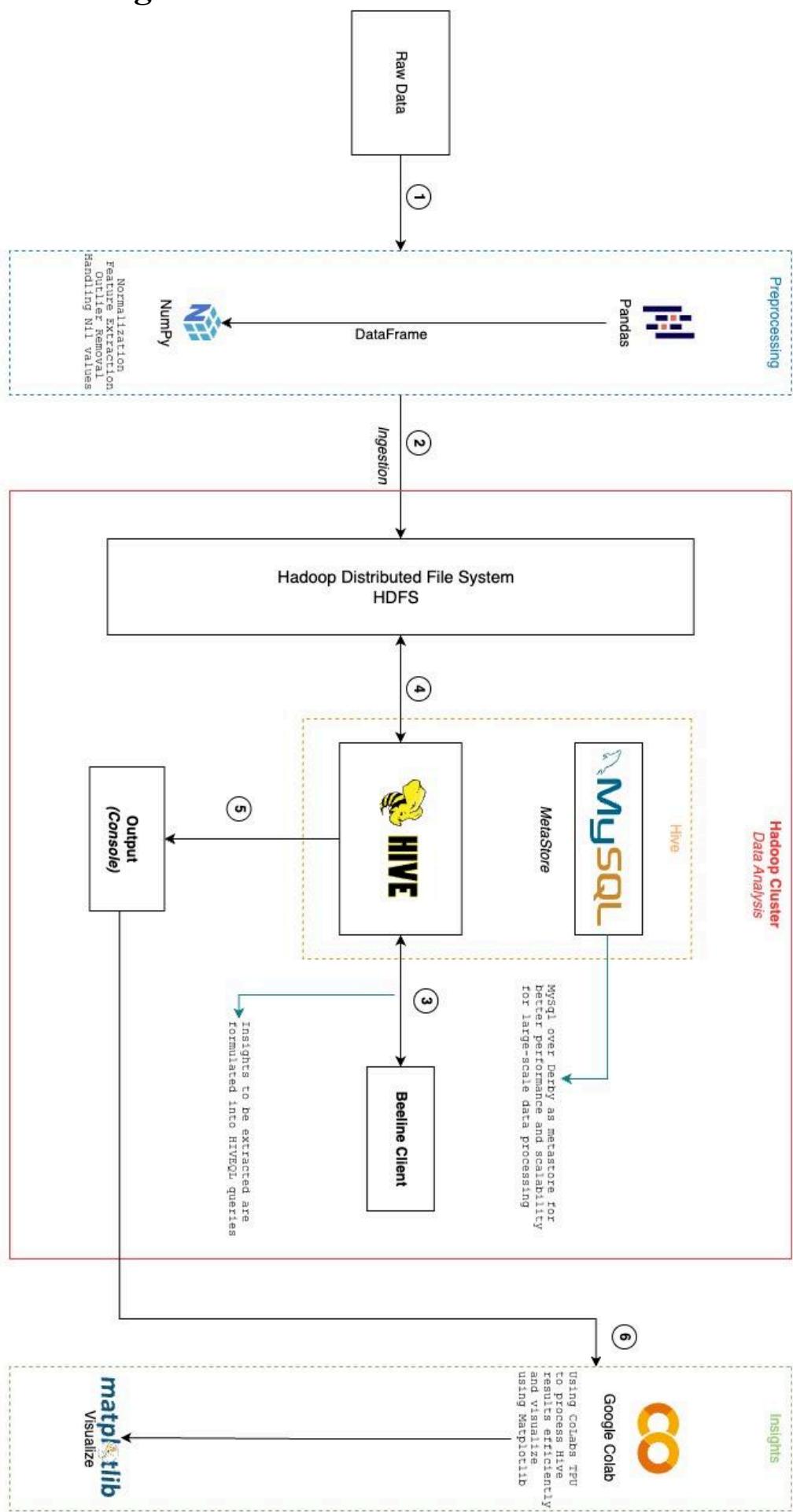
Hive's SQL-like query language, HiveQL, allows for complex data manipulations and aggregations, making it ideal for executing advanced SQL techniques like CTE's and window functions. These capabilities enable the extraction of meaningful insights from the dataset, such as identifying trends and correlations in commodity prices across different regions and seasons. Additionally, the use of Matplotlib complements Hive by transforming the query results into visualizations, providing a clear and intuitive representation of the data. This combination of tools ensures a comprehensive analysis, facilitating data-driven decision-making and enhancing the overall understanding of market dynamics.



## Objective

The main purpose of this project is to analyze the factors that influence the fluctuations in food prices in the Indian market with the help of big data analysis. By using advanced SQL techniques in Apache Hive, this study aims to uncover trends, correlations and patterns in raw material prices in various aspects such as geography and seasonality. The knowledge gained can inform stakeholders about market conditions, optimize pricing strategies, and improve supply chain efficiency. Ultimately, the project aims to contribute to a more stable and fair food market by conferring implementable data-controlled knowledge.

# Architecture Diagram



# Pathway of Taking the Project

The completion of this project takes place in a stepwise manner, to convert raw data into insightful information. The process is a series of meticulous steps aimed to provide an accurate and relevant final analysis. Here is a summary of the major steps included:

- **Dataset Acquisition**

A dataset containing Indian food price information was collected from Kaggle. This dataset acts as the basis for all later analysed data in terms of different commodities, geographical areas and time periods.

- **Preprocessing**

Before diving into analysis, we meticulously prepare our raw datasets using a combination of Pandas and NumPy. This crucial first step involves normalizing values to comparable scales, extracting meaningful features that drive insights, identifying and removing statistical outliers that could skew results, and developing intelligent strategies for handling missing data points. This ensures our foundation is solid before moving to heavyweight processing.

- **Data Ingestion**

Once cleaned, we migrate our datasets to Hadoop's distributed file system (HDFS). Unlike traditional storage solutions, HDFS spreads data across multiple nodes in the cluster, creating natural redundancy and enabling parallel processing capabilities. This architecture becomes particularly valuable as our dataset contains more than 160K records, providing both the reliability and performance needed for enterprise-scale analytics.

- **Data Analysis**

We leverage Hive's SQL-like interface to perform complex analytical operations within the Hadoop ecosystem. Our Hive Metastore runs on a dedicated MySQL backend, dramatically improving query performance and metadata management compared to Derby alternatives. We craft specialized HiveQL queries that help uncover hidden patterns and relationships within the data, focusing particularly on temporal trends and categorical correlations.

- **Visualize**

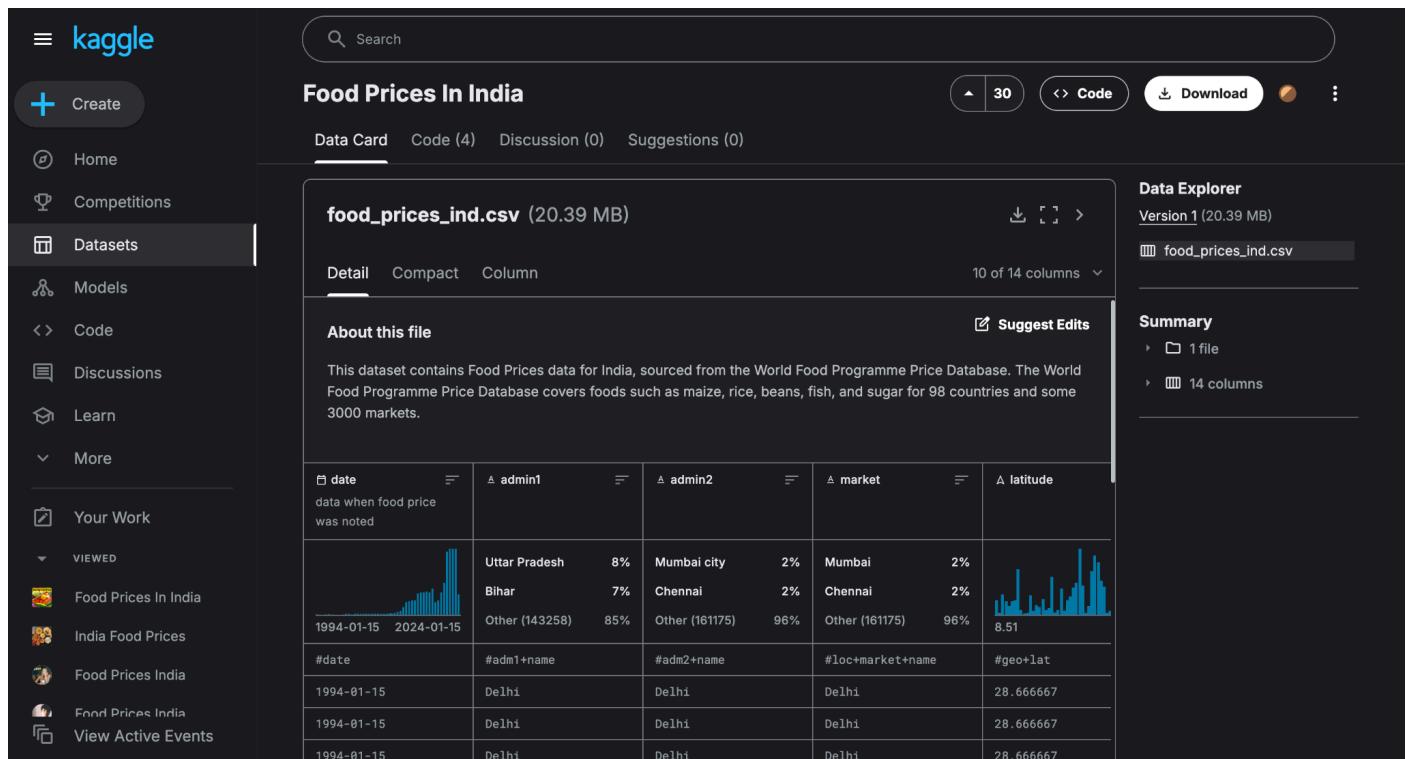
The final analytical outputs are transferred to Google Colab notebooks where we harness Matplotlib's extensive visualization capabilities and Colab's Tensor Processing Units. These visual assets transform raw numbers into compelling narratives that drive strategic decision-making.

# Dataset Description

The Indian Food Prices dataset is a comprehensive collection of food commodity price information, capturing over 1,60,000 records across multiple Indian states and cities, sourced from the World Food Programme Price Database. The World Food Programme Price Database covers foods such as maize, rice, beans, fish, and sugar for 98 countries and some 3000 markets.

The dataset includes detailed market-level data with attributes such as date, state, city, market location, geographical coordinates, commodity category, specific commodity, pricing unit, and price variations. It provides a granular view of food pricing dynamics, encompassing wholesale and retail prices across different markets, enabling in-depth analysis of regional price fluctuations, seasonal trends, and economic patterns in the Indian food commodity landscape. The data spans multiple years, offering a robust tool for understanding agricultural market behaviors, price correlations, and potential economic insights.

Source: <https://www.kaggle.com/datasets/abhinavshaw09/food-prices-in-india>



Screenshot from kaggle

## Dataset Attributes:

Field Name	Datatype	Description
date	STRING	The date on which the price was recorded.
state	STRING	The state in India where the price was observed.
city	STRING	The city within the state where the market is located.
market	STRING	The specific market or mandi where the commodity price was recorded.
latitude	FLOAT	The latitude coordinate of the market location.
longitude	FLOAT	The longitude coordinate of the market location.
category	STRING	The category of the commodity (e.g., grains, vegetables).
commodity	STRING	The specific commodity being priced (e.g., rice, wheat).
unit	STRING	The unit of measurement for the commodity price (e.g., kg, liter).
priceflag	STRING	An indicator of whether the price is estimated or actual.
currency	STRING	The currency in which the price is recorded. Mostly INR.
price	FLOAT	The recorded price of the commodity.
usdprice	FLOAT	The price converted to USD (if applicable).

*Below columns were generated during preprocessing of the dataset*

year	INTEGER	The year in which the price was recorded.
month	INTEGER	The month in which the price was recorded.
season	STRING	The season during which the price was recorded.

# Questions Overview

The analysis covers various dimensions such as geography, seasonality, and market conditions, providing a comprehensive understanding of the factors influencing food prices. We aim to uncover trends, correlations, and patterns that can inform strategic decision-making.

## Questions:

1. What is the average price of each commodity across different states, and how does it vary by season?
2. Which city has the highest average price for each commodity, and how does it compare to the national average price for that commodity?
3. Identify the top 5 commodities with the highest price volatility across different markets and states.
4. Which state has the most significant price difference between the highest and lowest priced commodities?
5. Which markets have the highest average price for each category, and how does it compare to the overall average price for that category?
6. Show only the Top 10 commodities with the highest price increase over the past year.
7. If a market has high prices for vegetables and fruits, does it also have high prices for cereals and tubers?
8. Rank commodities within each state based on their average price, and identify the top 3 most expensive commodities per state.
9. How do food prices in metropolitan cities compare to other markets during different seasons of the year?
10. Analyse the evolution of commodity prices over time, revealing year-to-year price dynamics and identifying significant market trends.
11. Identify the top 5 markets with the highest average price for each commodity, considering only the latest year.
12. Analyze the impact of seasonality on price volatility to calculate the coefficient of variation for each commodity.
13. Identify the cheapest and costliest state for each commodity for a given year? Take the year 2024 for example.
14. Which cities experienced consistent food price inflation over the years?
15. Identify the top 10 regions in India where food prices are the most volatile, by calculating the average price variance across nearby markets (within a 100 km radius).

## **Why we Selected these Questions:**

- To understand how prices fluctuate with seasons, helping stakeholders anticipate changes and plan inventory accordingly. This insight is crucial for managing supply and demand effectively.
- By analyzing price differences across states and cities, we can uncover regional disparities. This helps businesses tailor their strategies to local market conditions, ensuring competitive pricing.
- Examining how prices vary in different markets reveals the impact of local factors such as transportation costs and market demand. This knowledge allows for more precise pricing strategies.
- Focusing on individual commodities helps identify which items are most volatile or stable. This is vital for risk management and investment decisions.
- Analyzing price fluctuations over time uncovers patterns that can indicate market stability or instability. This insight is essential for long-term planning and forecasting.
- Understanding the relationships between different variables, such as location and price, helps identify underlying factors driving price changes. This can lead to more informed decision-making.

# Data Loading (Preparing Hive for the analysis)

```
-- Main Table
CREATE TABLE food_prices (
    date DATE,
    state STRING,
    city STRING,
    market STRING,
    latitude DOUBLE,
    longitude DOUBLE,
    category STRING,
    commodity STRING,
    unit STRING,
    price DOUBLE,
    year INT,
    month INT,
    season STRING
) PARTITIONED BY (category STRING) CLUSTERED BY (state) INTO 10 BUCKETS STORED AS
PARQUET;

-- External table for data loading from hdfs
CREATE TABLE staging_food_prices (
    `date` DATE,
    state STRING,
    city STRING,
    market STRING,
    latitude DOUBLE,
    longitude DOUBLE,
    category STRING,
    commodity STRING,
    unit STRING,
    priceflag STRING,
    pricetype STRING,
    currency STRING,
    price DOUBLE,
    usdprice DOUBLE,
    year INT,
    month INT,
    season STRING
) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',';

LOAD DATA INPATH '/data/preprocessed_india_food_prices.csv' INTO TABLE
staging_food_prices;

-- Skip the csv header while selecting
ALTER TABLE staging_food_prices
SET TBLPROPERTIES ("skip.header.line.count" = "1");

-- Enable dynamic partitioning
SET hive.exec.dynamic.partition=true;
SET hive.exec.dynamic.partition.mode=nonstrict;

-- Insert from staging to partitioned table
INSERT OVERWRITE TABLE food_prices PARTITION(category)
SELECT
```

```

`date`,
state,
city,
market,
latitude,
longitude,
commodity,
unit,
price,
year,
month,
season,
category
FROM staging_food_prices;

SELECT COUNT(*) FROM food_prices;

```

```

SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/hive/lib/log4j-slf4j-impl-2.6.2.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/hadoop-2.7.4/share/hadoop/common/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Connecting to jdbc:hive2://localhost:10000
Connected to: Apache Hive (version 2.3.2)
Driver: Hive JDBC (version 2.3.2)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 2.3.2 by Apache Hive
0: jdbc:hive2://localhost:10000> use case_study;
No rows affected (0.197 seconds)
0: jdbc:hive2://localhost:10000> CREATE TABLE food_prices ( `date` DATE, state STRING, city STRING, market STRING, latitude DOUBLE, longitude DOUBLE, commodity STRING, unit STRING, price DOUBLE, year INT, month INT, season STRING ) PARTITIONED BY (state) INTO 10 BUCKETS STORED AS PARQUET;
No rows affected (0.411 seconds)
0: jdbc:hive2://localhost:10000> SET hive.exec.dynamic.partition=true;
No rows affected (0.037 seconds)
0: jdbc:hive2://localhost:10000> SET hive.exec.dynamic.partition.mode=nonstrict;
No rows affected (0.019 seconds)
0: jdbc:hive2://localhost:10000> CREATE TABLE staging_food_prices ( `date` DATE, state STRING, city STRING, market STRING, latitude DOUBLE, longitude DOUBLE, category STRING, commodity STRING, unit STRING, priceflag STRING, pricetype STRING, currency STRING, price DOUBLE, usdprice DOUBLE, year INT, month INT, season STRING ) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',';
No rows affected (0.697 seconds)
0: jdbc:hive2://localhost:10000> LOAD DATA INPATH '/data/preprocessed_india_food_prices.csv' INTO TABLE staging_food_prices;
No rows affected (1.526 seconds)
0: jdbc:hive2://localhost:10000> ALTER TABLE staging_food_prices SET TBLPROPERTIES ("skip.header.line.count"="1");
No rows affected (0.353 seconds)
0: jdbc:hive2://localhost:10000> INSERT OVERWRITE TABLE food_prices PARTITION(category) SELECT `date`, state, city, market, latitude, longitude, commodity, unit, price, year, month, season, category FROM staging_food_prices;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
No rows affected (25.407 seconds)
0: jdbc:hive2://localhost:10000> Display all 560 possibilities? (y or n)
0: jdbc:hive2://localhost:10000> DROP TABLE staging_food_prices;
No rows affected (1.593 seconds)
0: jdbc:hive2://localhost:10000> SELECT COUNT(*) FROM food_prices;
+-----+
| _c0 |
+-----+
| 166096 |
+-----+
1 row selected (1.238 seconds)
0: jdbc:hive2://localhost:10000> ■

```

# Data Analysis

Q1) What is the average price of each commodity across different states, and how does it vary by season?

```
SELECT
    commodity,
    state,
    season,
    AVG(price) AS avg_price
FROM
    food_prices
GROUP BY
    commodity, state, season
ORDER BY
    commodity, state, season;
```

commodity	state	season	avg_price
Chickpeas	Bihar	Autumn	43.699803921568616
Chickpeas	Bihar	Monsoon	40.633698652173985
Chickpeas	Bihar	Summer	41.48464755319149
Chickpeas	Bihar	Winter	41.48348837209302
Chickpeas	Delhi	Autumn	54.0311363336336344
Chickpeas	Delhi	Monsoon	58.1788
Chickpeas	Delhi	Summer	42.98411764795883
Chickpeas	Delhi	Winter	46.83926829268295
Chickpeas	Maharashtra	Autumn	52.022883333333331
Chickpeas	Maharashtra	Monsoon	48.89627456980393
Chickpeas	Maharashtra	Summer	46.523333333333326
Chickpeas	Maharashtra	Winter	55.50717948717948
Chickpeas	Tamil Nadu	Autumn	50.016000000000005
Chickpeas	Tamil Nadu	Monsoon	44.983673464938776
Chickpeas	Tamil Nadu	Summer	47.099999999999994
Chickpeas	Tamil Nadu	Winter	44.425106382978726
Ghee (vanaspati)	Andaman And Nicobar	Autumn	162.618125
Ghee (vanaspati)	Andaman And Nicobar	Monsoon	154.72625
Ghee (vanaspati)	Andaman And Nicobar	Summer	157.32176470588234
Ghee (vanaspati)	Andaman And Nicobar	Winter	167.016
Ghee (vanaspati)	Andhra Pradesh	Autumn	103.05545454545454
Ghee (vanaspati)	Andhra Pradesh	Monsoon	103.05545454545453
Ghee (vanaspati)	Andhra Pradesh	Summer	106.55314285714286
Ghee (vanaspati)	Andhra Pradesh	Winter	109.11545454545454
Ghee (vanaspati)	Assam	Autumn	107.408633636363635
Ghee (vanaspati)	Assam	Monsoon	105.86571428571429
Ghee (vanaspati)	Assam	Summer	98.86468998998992
Ghee (vanaspati)	Assam	Winter	98.86468998998992
Ghee (vanaspati)	Bihar	Autumn	120.54246262626265
Ghee (vanaspati)	Bihar	Monsoon	122.38546262626265
Ghee (vanaspati)	Bihar	Summer	126.0941221374944
Ghee (vanaspati)	Bihar	Winter	121.7838882352941
Ghee (vanaspati)	Chandigarh	Autumn	98.146479582352941
Ghee (vanaspati)	Chandigarh	Monsoon	109.783333333333333
Ghee (vanaspati)	Chandigarh	Summer	102.55268649656522
Ghee (vanaspati)	Chandigarh	Winter	101.34235294171647
Ghee (vanaspati)	Chhattisgarh	Autumn	116.69718750000001
Ghee (vanaspati)	Chhattisgarh	Monsoon	129.2775
Ghee (vanaspati)	Chhattisgarh	Summer	133.7174874874874
Ghee (vanaspati)	Chhattisgarh	Winter	122.94769000000001
Ghee (vanaspati)	Delhi	Autumn	111.41571428571429
Ghee (vanaspati)	Delhi	Monsoon	115.699
Ghee (vanaspati)	Delhi	Summer	117.37857142857142
Ghee (vanaspati)	Delhi	Winter	106.21416666666669
Ghee (vanaspati)	Goa	Autumn	125.914375
Ghee (vanaspati)	Goa	Monsoon	142.6423876923075
Ghee (vanaspati)	Goa	Summer	132.63277777777776
Ghee (vanaspati)	Goa	Winter	131.3790909090909
Ghee (vanaspati)	Gujarat	Autumn	118.88227272727272
Ghee (vanaspati)	Gujarat	Monsoon	106.2840277777778
Ghee (vanaspati)	Gujarat	Summer	110.80283582889552
Ghee (vanaspati)	Gujarat	Winter	110.816666666666668
Ghee (vanaspati)	Haryana	Autumn	106.0271186440678
Ghee (vanaspati)	Haryana	Monsoon	101.3058241758242
Ghee (vanaspati)	Haryana	Summer	104.7451612983257
Ghee (vanaspati)	Haryana	Winter	103.01914634146338
Ghee (vanaspati)	Himachal Pradesh	Autumn	122.0197951836736
Ghee (vanaspati)	Himachal Pradesh	Monsoon	122.96983050847462

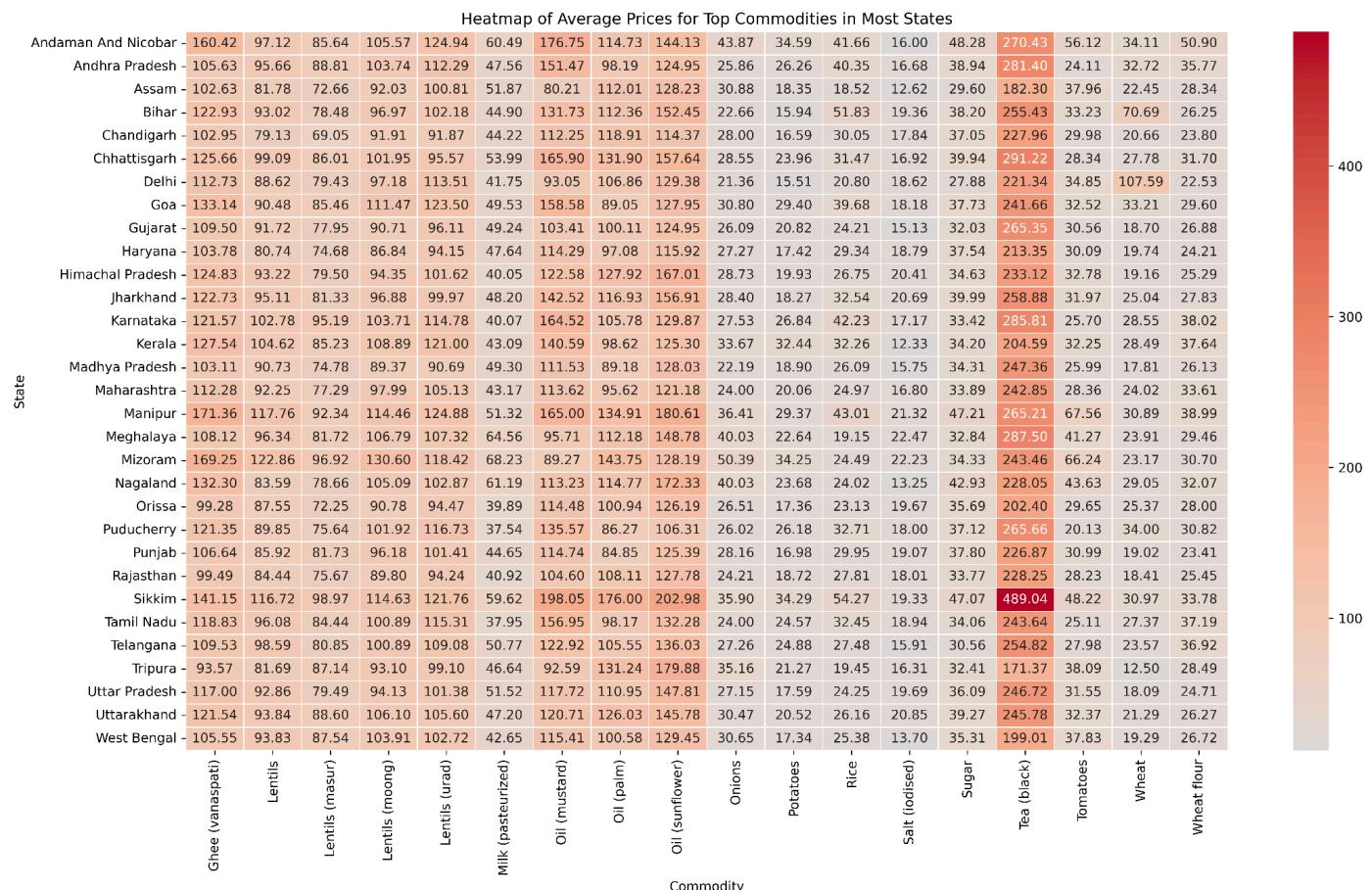
## Relevance

Understanding average commodity prices across states and seasons reveals crucial regional and temporal price patterns. This helps identify which states consistently experience higher food costs and when prices typically peak or dip throughout the year. Such data helps both

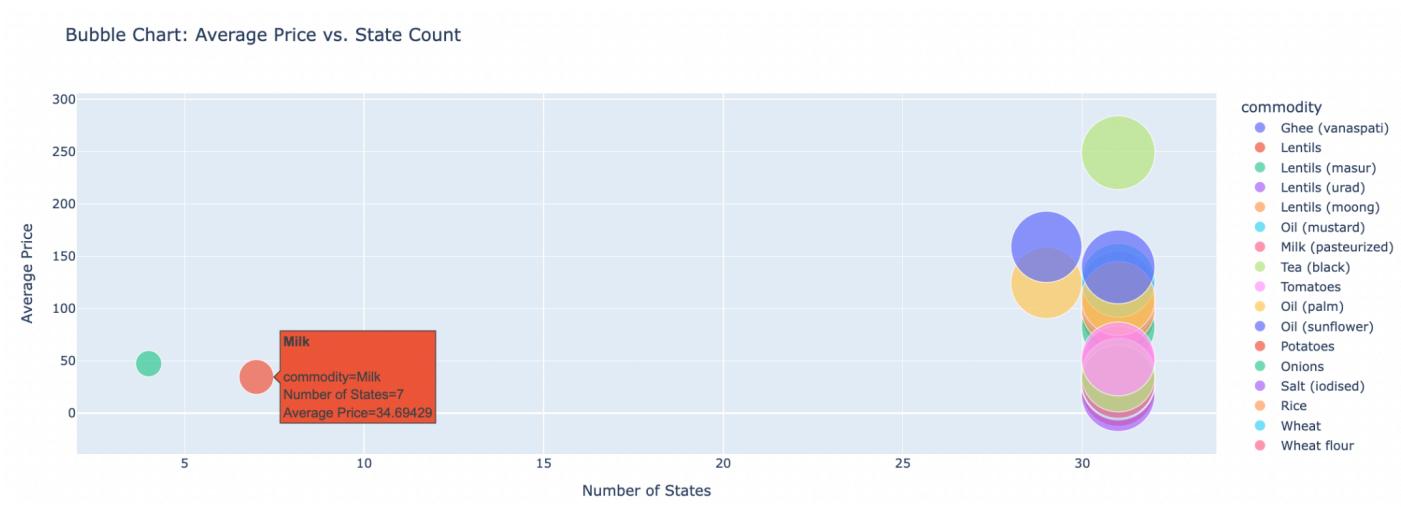
consumers and policymakers manage their budgets and determine where and when to deploy food security measures.

## Visualization

Heatmap showing the average prices of the commodities that are present in most states.



The bubble chart shows the relationship between the number of states a commodity is present in and its average price. Larger bubbles indicate commodities that are available in more states, while the y-axis represents their average price.



Q2) Which city has the highest average price for each commodity, and how does it compare to the national average price for that commodity?

```
WITH national_avg AS (
    SELECT
        commodity,
        AVG(price) AS national_avg_price
    FROM
        food_prices
    GROUP BY
        commodity
)
SELECT
    a.city,
    a.commodity,
    AVG(a.price) AS city_avg_price,
    b.national_avg_price
FROM
    food_prices a
JOIN
    national_avg b ON a.commodity = b.commodity
GROUP BY
    a.city, a.commodity, b.national_avg_price
ORDER BY
    a.commodity, city_avg_price DESC;
```

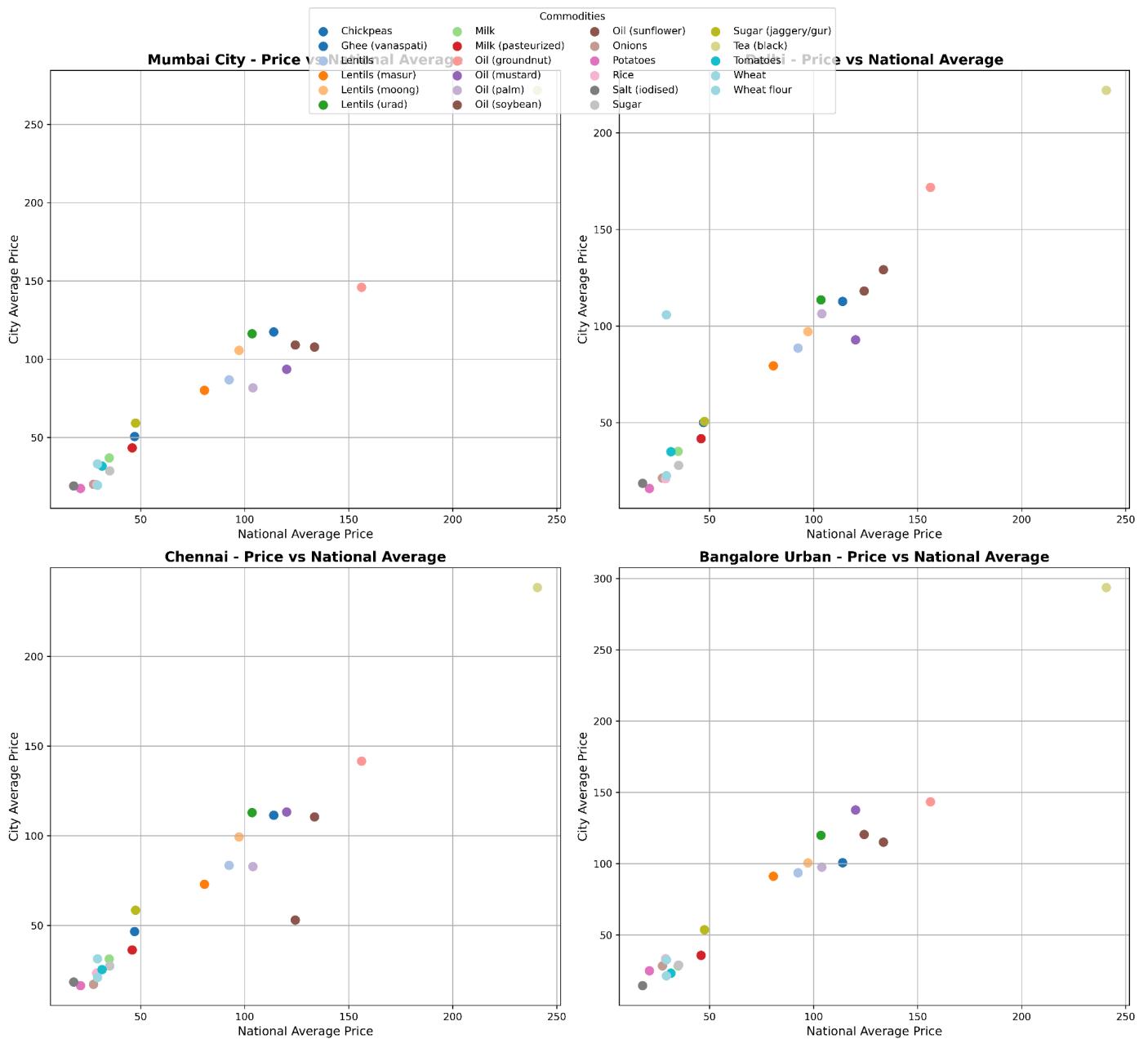
a.city	a.commodity	city_avg_price	b.national_avg_price
Mumbai City	Chickpeas	50.698500000000014	47.07922877922878
Delhi	Chickpeas	50.074319444445	47.07922877922878
Chennai	Chickpeas	46.56943956843957	47.07922877922878
Patna	Chickpeas	41.3798374331551	47.07922877922878
Mysore	Ghee (vanaspati)	207.25678571428574	113.96868593289889
Kohima	Ghee (vanaspati)	178.84714285714287	113.96868593289889
Wayanad	Ghee (vanaspati)	178.667777777778	113.96868593289889
West Imphal	Ghee (vanaspati)	178.1375	113.96868593289889
Aizawl	Ghee (vanaspati)	168.62	113.96868593289889
Valsad	Ghee (vanaspati)	161.21666666666667	113.96868593289889
Andaman Islands	Ghee (vanaspati)	159.01962962962966	113.96868593289889
Saharanpur	Ghee (vanaspati)	158.8345	113.96868593289889
Thrissur	Ghee (vanaspati)	156.56222222222223	113.96868593289889
Palakkad	Ghee (vanaspati)	153.61862868965513	113.96868593289889
Jhabua	Ghee (vanaspati)	153.26769238769233	113.96868593289889
Vadodara	Ghee (vanaspati)	152.6717391304348	113.96868593289889
Munger	Ghee (vanaspati)	152.35071428571428	113.96868593289889
Purba Singhbhum	Ghee (vanaspati)	151.935652173913	113.96868593289889
Saran	Ghee (vanaspati)	151.21384615384613	113.96868593289889
Araria	Ghee (vanaspati)	151.04272727272726	113.96868593289889
The Dangs	Ghee (vanaspati)	150.6221052631579	113.96868593289889
Jaintia Hills	Ghee (vanaspati)	149.82424242424242	113.96868593289889
Katihar	Ghee (vanaspati)	149.08833333333334	113.96868593289889
Meerut	Ghee (vanaspati)	149.20033233333334	113.96868593289889
Reigarh	Ghee (vanaspati)	149.80789732694212	113.96868593289889
Bijapur	Ghee (vanaspati)	149.8694117647859	113.96868593289889
Amravati	Ghee (vanaspati)	148.85111111111115	113.96868593289889
Lohardaga	Ghee (vanaspati)	148.47466666666665	113.96868593289889
Cuddalore	Ghee (vanaspati)	148.43468333333333	113.96868593289889
Mirzapur	Ghee (vanaspati)	148.16511111111112	113.96868593289889
Solan	Ghee (vanaspati)	146.00150000000002	113.96868593289889
Chamba	Ghee (vanaspati)	145.72812499999998	113.96868593289889
Darrang	Ghee (vanaspati)	145.24166666666667	113.96868593289889
Aligarh	Ghee (vanaspati)	144.42736842185262	113.96868593289889
Kangra	Ghee (vanaspati)	144.49166666666667	113.96868593289889
Vellore	Ghee (vanaspati)	144.35387692387694	113.96868593289889
Hoshangabad	Ghee (vanaspati)	144.0776923876923	113.96868593289889
Bareilly	Ghee (vanaspati)	143.87478268689565	113.96868593289889
Naini Tal	Ghee (vanaspati)	143.73889523889524	113.96868593289889
Una	Ghee (vanaspati)	143.3441176478582	113.96868593289889
Madhubani	Ghee (vanaspati)	143.32500000000002	113.96868593289889
Haridwar	Ghee (vanaspati)	142.33318181818183	113.96868593289889

## Relevance

Identifying cities with the highest costs for each product highlights urban areas with the most severe cost-of-living constraints. Comparing these prices to national averages reveals how much more inhabitants in these areas spend for necessities, which is useful for understanding urban poverty and designing targeted subsidies or market interventions.

## Visualization

The visualization shows a price comparison using a scatter plot for 4 major Indian cities (Mumbai City, Delhi, Chennai and Bengaluru Urban). Results from the above query were processed in Python to draw the relevant visualization.



Q3) Identify the top 5 commodities with the highest price volatility across different markets and states.

```
SELECT
    commodity,
    STDDEV(price) AS price_volatility
FROM
    food_prices
GROUP BY
    commodity
ORDER BY
    price_volatility DESC
LIMIT 5;
```

```
0: jdbc:hive2://localhost:10000> SELECT commodity, STDDEV(price) AS price_volatility FROM food_prices GROUP BY commodity ORDER BY price_volatility DESC LIMIT 5;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
+-----+-----+
| commodity | price_volatility |
+-----+-----+
| Tea (black) | 66.2896643622414 |
| Wheat | 65.61679154387093 |
| Oil (mustard) | 48.583769518282615 |
| Oil (sunflower) | 38.58468014708176 |
| Ghee (vanaspati) | 35.30627425568045 |
+-----+
5 rows selected (4.955 seconds)
0: jdbc:hive2://localhost:10000> ■
```

## Relevance

Price volatility has a direct impact on food security and household budgets. By finding the top five commodities with the highest volatility, we can determine which foods cause the most financial worry for households. This information is crucial for government agencies focused on price stabilisation and organisations working on food security programs.

Q4) Which state has the most significant price difference between the highest and lowest priced commodities?

```
WITH state_price_range AS (
    SELECT
        state,
        MAX(price) AS max_price,
        MIN(price) AS min_price
    FROM
        food_prices
    GROUP BY
        state
)
SELECT
    state,
    (max_price - min_price) AS price_difference
FROM
    state_price_range
ORDER BY
    price_difference DESC
LIMIT 1;
```

```
0: jdbc:hive2://localhost:10000> WITH state_price_range AS ( SELECT state, MAX(price) AS max_price, MIN(price) AS min_price FROM food_prices GROUP BY state ) SELECT state, (max_price - min_price) AS price_difference FROM state_price_range ORDER BY price_difference DESC LIMIT 1;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
+-----+-----+
| state | price_difference |
+-----+-----+
| Bihar | 740.19          |
+-----+-----+
1 row selected (4.816 seconds)
0: jdbc:hive2://localhost:10000> ■
```

## Relevance

The state with the largest price difference between commodities likely experiences significant economic disparity or logistical challenges. This analysis highlights where certain foods might be luxury items while others remain affordable, pointing to possible supply chain inefficiencies or regional production specialization that affects accessibility of diverse foods.

Q5) Which markets have the highest average price for each category, and how does it compare to the overall average price for that category?

```
WITH category_avg AS (
    SELECT
        category,
        AVG(price) AS overall_avg_price
    FROM
        food_prices
    GROUP BY
        category
)
SELECT
    a.market,
    a.category,
    AVG(a.price) AS market_avg_price,
    b.overall_avg_price
FROM
    food_prices a
JOIN
    category_avg b ON a.category = b.category
GROUP BY
    a.market, a.category, b.overall_avg_price
ORDER BY
    a.category, market_avg_price DESC;
```

```
0: jdbc:hive2://localhost:10000> WITH category_avg AS ( SELECT category, AVG(price) AS overall_avg_price FROM food_prices GROUP BY category ) SELECT a.market, a.category, AVG(a.price) AS market_avg_price, b.overall_avg_price FROM food_prices a JOIN category_avg b ON a.category = b.category GROUP BY a.market, a.category, b.overall_avg_price ORDER BY a.category, market_avg_price DESC;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/hive/lib/log4j-slf4j-impl-2.6.2.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/hadoop/common/lib/slf4j-log4j1-2.1.7.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Execution log at: /tmp/root/root_20250322135656_ebad5668-dc35-44ab-8c69-93d8147ab51d.log
2025-03-22 13:57:05 Starting to launch local task to process map join; maximum memory = 477626368
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Execution log at: /tmp/root/root_20250322135656_ebad5668-dc35-44ab-8c69-93d8147ab51d.log
2025-03-22 13:57:05 Starting to launch local task to process map join; maximum memory = 477626368
SLF4J: Failed to load class "org.slf4j.impl.StaticLoggerBinder".
SLF4J: Defaulting to no-operation (NOP) logger implementation
SLF4J: See http://www.slf4j.org/codes.html#staticLoggerBinder for further details.
2025-03-22 13:57:07 Dump the side-table for tag: @ with group count: 6 into file: file:/tmp/root/cf169404-99eb-45eb-b479-2d9ae48f96f4/hive_2025-03-22_13-56-56_721_6876095130143795652-1-local-10007/Ha
shTable-Stage-2/MapJoin-mapfile20--.hashtable
2025-03-22 13:57:08 Uploaded 1 File to: file:/tmp/root/cf169404-99eb-45eb-b479-2d9ae48f96f4/hive_2025-03-22_13-56-56_721_6876095130143795652-1-local-10007/HashTable-Stage-2/MapJoin-mapfile20--.hashta
ble (3620529 bytes)
2025-03-22 13:57:08 End of local task; Time Taken: 3.145 sec.
```

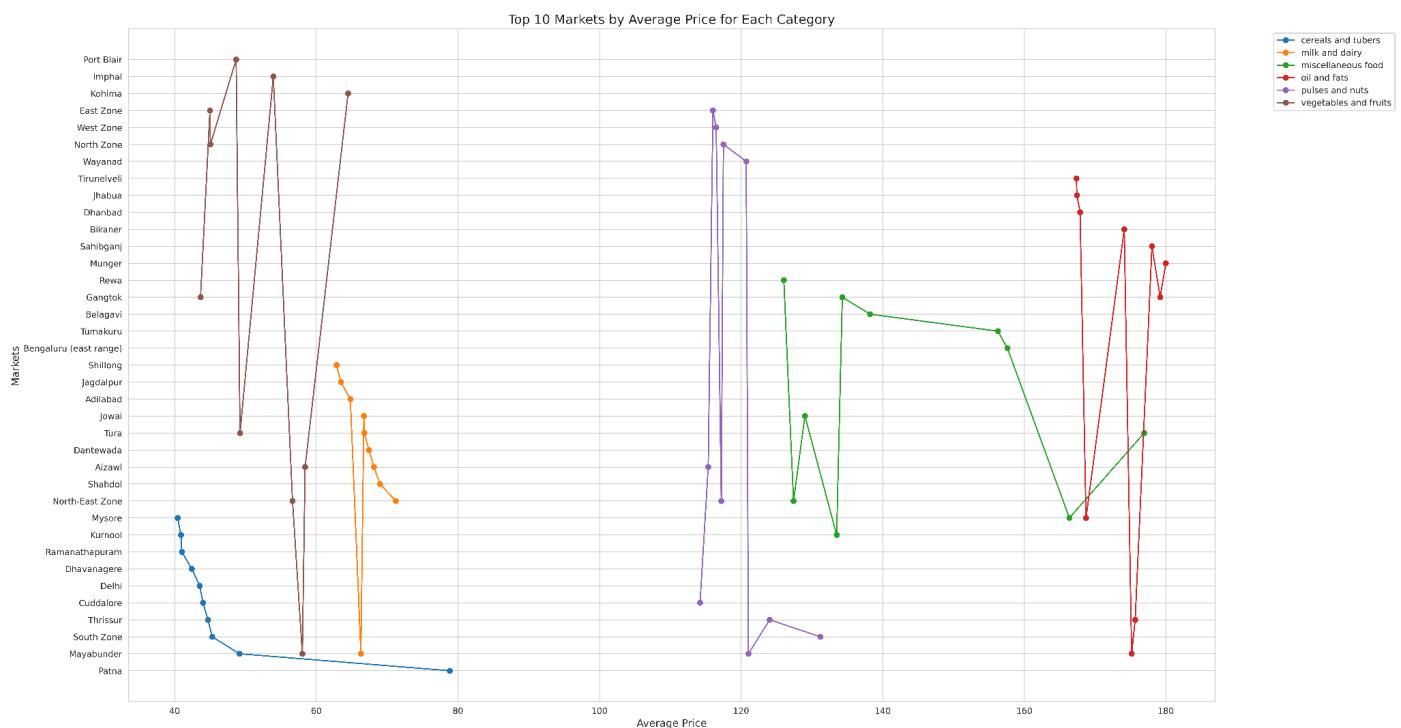
a.market	a.category	market_avg_price	b.overall_avg_price
Patna	cereals and tubers	78.85346652267819	27.28121837279347
Maybunder	cereals and tubers	49.14235789757	27.28121837279347
South Zone	cereals and tubers	45.298260849865226	27.28121837279347
Thrissur	cereals and tubers	44.699907467407416	27.28121837279347
Cuddalore	cereals and tubers	43.995946594859485	27.28121837279347
Delhi	cereals and tubers	43.537741935483886	27.28121837279347
Dhavanagere	cereals and tubers	42.4668885186383	27.28121837279347
Ramanathapuram	cereals and tubers	41.807183399891554	27.28121837279347
Kurnool	cereals and tubers	40.875769239876924	27.28121837279347
Mysore	cereals and tubers	40.42229357798165	27.28121837279347
Tirupathi	cereals and tubers	39.64473214285715	27.28121837279347
Tirunelveli	cereals and tubers	39.583333333333336	27.28121837279347
Adilabad	cereals and tubers	39.48296296296297	27.28121837279347
Vellore	cereals and tubers	39.244736842105276	27.28121837279347
Maharashtra_Pune	cereals and tubers	39.210927835865154	27.28121837279347
Palakkad	cereals and tubers	38.969667796618176	27.28121837279347
Wayanad	cereals and tubers	38.57773584985661	27.28121837279347
Selambra	cereals and tubers	38.254125	27.28121837279347
Mangalore	cereals and tubers	38.007680000000001	27.28121837279347
Port Blair	cereals and tubers	37.62918421052631	27.28121837279347
Gangtok	cereals and tubers	37.47848837289303	27.28121837279347
Coimbatore	cereals and tubers	37.40037837037836	27.28121837279347
Bellary	cereals and tubers	37.211914893617035	27.28121837279347
Dharmapuri	cereals and tubers	37.114186046511634	27.28121837279347
Kalaburagi	cereals and tubers	37.06673469387756	27.28121837279347
Tumakuru	cereals and tubers	37.05448979591838	27.28121837279347
Mangaon-raigad	cereals and tubers	36.960243962439025	27.28121837279347
North-East Zone	cereals and tubers	36.77045454545454	27.28121837279347
T.Puram	cereals and tubers	36.4801637837804	27.28121837279347
Imphal	cereals and tubers	36.49445568400785	27.28121837279347
Pettewada	cereals and tubers	36.22393674800755	27.28121837279347
Vijayapur	cereals and tubers	35.9959803921568	27.28121837279347
Shivamogga	cereals and tubers	35.73820408163264	27.28121837279347
Kozhikode	cereals and tubers	35.6900784213754	27.28121837279347
Bellavi	cereals and tubers	35.49911764798581	27.28121837279347
Ajmer	cereals and tubers	35.46107692387692	27.28121837279347
Bengaluru (east range)	cereals and tubers	35.4158762884598	27.28121837279347
Kohima	cereals and tubers	35.40796477419355	27.28121837279347
Bikaner	cereals and tubers	34.8436170212766	27.28121837279347
Karimnagar	cereals and tubers	34.78384347826987	27.28121837279347
Warangal	cereals and tubers	34.450101010101015	27.28121837279347
Tura	cereals and tubers	33.677592592592596	27.28121837279347

## Relevance

Markets with the highest average prices for food categories often indicate areas where residents face greater food insecurity risks. Comparing these to overall category averages reveals the magnitude of the cost burden in these locations. This information helps identify communities that might need targeted support through public distribution systems or improved market access.

## Visualization

This line chart visualizes the top 10 markets for each product category, displaying their average prices across different regions of India.



Q6) Show only the Top 10 commodities with the highest price increase over the past year.

```
WITH yearly_avg AS (
    SELECT
        commodity,
        year,
        AVG(price) AS avg_price
    FROM
        food_prices
    GROUP BY
        commodity, year
),
price_increase AS (
    SELECT
        a.commodity,
        (a.avg_price - b.avg_price) AS price_diff
    FROM
        yearly_avg a
    JOIN
        yearly_avg b ON a.commodity = b.commodity AND a.year = b.year + 1
)
SELECT
    commodity,
    price_diff
FROM
    price_increase
ORDER BY
    price_diff DESC
LIMIT 10;
```

```
0: jdbc:hive2://localhost:10000> WITH yearly_avg AS ( SELECT commodity, year, AVG(price) AS avg_price FROM food_prices GROUP BY commodity, year ), price_increase AS ( SELECT a.commodity, (a.avg_price - b.avg_price) AS price_diff FROM yearly_avg a JOIN yearly_avg b ON a.commodity = b.commodity AND a.year = b.year + 1 ) SELECT commodity, price_diff FROM price_increase ORDER BY price_diff DESC LIMIT 10;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/hive/lib/log4j-slf4j-impl-2.6.2.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/hadoop-2.7.4/share/hadoop/common/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Execution log at: /tmp/root/root_20250322140316_7fa51c5-5978-4073-aa94-7af03ad49944.log
2025-03-22 14:03:26     Starting to launch local task to process map join;           maximum memory = 477626368
2025-03-22 14:03:26     Dump the side-table for tag: 1 with group count: 379 into file: file:/tmp/root/cf169404-99eb-45eb-b479-2d9ae48f96f4/hive_2025-03-22_14-03-16_311_5219764060382328263-1/-local-10007/
HashTable-Stage-5/MapJoin-mapfile31--.hashtable
2025-03-22 14:03:26     Uploaded 1 File to: file:/tmp/root/cf169404-99eb-45eb-b479-2d9ae48f96f4/hive_2025-03-22_14-03-16_311_5219764060382328263-1/-local-10007/HashTable-Stage-5/MapJoin-mapfile31--.hashtable (15441 bytes)
2025-03-22 14:03:26     End of local task; Time Taken: 0.398 sec.
+-----+-----+
| commodity | price_diff |
+-----+-----+
| Wheat      | 74.25138832499814 |
| Wheat      | 69.45079483458997 |
| Tea (black) | 54.04354722338554 |
| Oil (mustard)| 45.09187584745192 |
| Oil (sunflower) | 44.20934129464681 |
| Rice       | 42.78775641825642 |
| Oil (soybean) | 41.03809827463114 |
| Chickpeas  | 40.6633227344992 |
| Ghee (vanaspati) | 35.16478043477676 |
| Lentils    | 33.01498920547865 |
+-----+-----+
10 rows selected (13.254 seconds)
0: jdbc:hive2://localhost:10000>
```

## Relevance

Tracking commodities with the highest price increases spotlights the emerging affordability crisis. These rapid price jumps can signal supply shortages, distribution problems, or changing consumption patterns. Policymakers and market regulators need this information to monitor inflation triggers and implement timely interventions before staple foods become unaffordable.

Q7) If a market has high prices for vegetables and fruits, does it also have high prices for cereals and tubers?

```
WITH category_avg AS (
    SELECT market, category, AVG(price) AS avg_price
    FROM food_prices
    WHERE category IN ('vegetables and fruits', 'cereals and tubers')
    GROUP BY market, category
)
SELECT market,
    MAX(CASE WHEN category IN ('vegetables and fruits') THEN avg_price ELSE 0 END) AS veg_fruit_price,
    MAX(CASE WHEN category IN ('cereals and tubers') THEN avg_price ELSE 0 END) AS cereal_tuber_price
FROM category_avg
GROUP BY market
HAVING veg_fruit_price > 0 AND cereal_tuber_price > 0
ORDER BY veg_fruit_price DESC, cereal_tuber_price DESC;
```

market	veg_fruit_price	cereal_tuber_price
Kohima	64.48636363636363	35.40709677419385
Aizawl	58.4173964497804135	28.4384615384618154
Mayabunder	58.02531914893619	49.1423157894737
North-East Zone	56.64	36.77045454545454
Imphal	53.93150000000001	36.33415384615385
Tura	49.233861224489795	33.677592592592596
Port Blair	48.67621359223381	37.62918421852631
North Zone	45.05454545454545	29.44818181818182
East Zone	44.975	31.33285714285714
Gangtok	43.652549819407856	37.178574128574134
Jamshedpur	41.79399999999992	30.4626229508194
Jowai	48.0212121212156	28.9426229508194
West Zone	39.75164646464645	30.02010201020102
Dibrugarh	39.68563451776651	25.991705882352942
Survept	38.17545454545455	32.83356521739131
Shillong	38.12506997560974	29.893770936223255
Dharmsala	37.85422222222222	24.6935826895522744
Raiganj	37.414423876923886	25.469348837289303
Thrisur	37.080192307692315	44.4999874074747416
Agartala	36.83878504672898	21.493638253635825
South Zone	36.6935353535353564	45.2982686869565226
Ramanathapuram	36.4565306122449	41.8071830988591554
Mangano-raigad	35.889767441869464	36.696243987439825
Sahibganj	35.837619047619036	24.25642857342857
Malda	35.587254981968788	24.998437499999998
Guwahati	35.4813888888888896	28.078164986937387
Wayanad	35.224888888888884	38.57773584905661
Tiruvananthapuram	35.1824193548387	18.789526184538644
Purulia	34.857391304347836	22.717619847619846
Haldwani	34.75887096774194	27.77558847457628
Gorakhpur	34.556522641509435	23.165789473684207
Allahabad	34.15063492063492	27.689669421487594
Warangal	34.06282608695653	34.450101010101015
Bikaner	33.697600000000001	34.8436170212766
Darbhanga	33.6605306122449	24.579157894736845
T.Puram	33.66082089552238	36.40783763703784
Balasore	33.62155172413794	27.286039403940393
Kharagpur	33.59275862068966	36.095315315315308
Gumla	33.58441176470588	28.739227988505747
Saharsa	33.52680851863829	28.79976923976922
Haridwar	33.4881428571429	26.48786781754387
Beripada	33.480439999999996	26.70079000000000002
Siliguri	33.31627457457457	14.27
Ernakulam	33.30257675675677	33.43087434554975
Chittewadu	33.19790697474418	36.229236708866758
Solan	33.08734285714287	36.095243990243993
Bokaro	32.9449717948718	25.782437793181345
Jeyapore	32.62423728813561	28.89757441869464
Dhanbad	32.354615384615386	28.25348837289382
Surat	32.349999999999994	27.936395348372905
Palakkad	32.29886955651739	38.969867796610176
Panaji	31.753825641825648	32.7913899192280857
Araria	31.746800000000002	27.994363797448362
Berhampur	31.62954545454545	28.776581196581194
Kolkata	31.5481208805369134	18.408392857142854
Rampurhat	31.480377358490568	30.56127906976746

## Relevance

The correlation between prices across different food groups reveals important market dynamics. If a market shows consistently high prices across categories, this suggests underlying structural issues like transportation bottlenecks or higher local operating costs, rather than commodity-specific problems. This helps distinguish between localized supply chain issues and broader regional economic challenges.

Q8) Rank commodities within each state based on their average price, and identify the top 3 most expensive commodities per state.

```
WITH ranked_prices AS (
    SELECT
        state,
        commodity,
        AVG(price) AS avg_price,
        RANK() OVER (PARTITION BY state ORDER BY AVG(price) DESC) AS price_rank
    FROM food_prices
    GROUP BY state, commodity
)
SELECT
    state,
    commodity,
    avg_price,
    price_rank
FROM ranked_prices
WHERE price_rank <= 3
ORDER BY state, price_rank;
```

state	commodity	avg_price	price_rank
Andaman And Nicobar	Tea (black)	271.0917924528302	1
Andaman And Nicobar	Oil (mustard)	175.07275	2
Andaman And Nicobar	Ghee (vanaspati)	159.01962962962966	3
Andhra Pradesh	Tea (black)	281.17457427118437	1
Andhra Pradesh	Oil (mustard)	151.1030106145251	2
Andhra Pradesh	Oil (groundnut)	142.51578707070712	3
Assam	Oil (groundnut)	194.61777777777777	1
Assam	Tea (black)	181.744695655217393	2
Assam	Oil (sunflower)	128.18233766233767	3
Bihar	Tea (black)	255.38103448275876	1
Bihar	Oil (groundnut)	188.4385607476635	2
Bihar	Oil (sunflower)	162.55872817353673	3
Chandigarh	Tea (black)	227.92811965811964	1
Chandigarh	Oil (groundnut)	139.12876923876924	2
Chandigarh	Oil (palm)	118.8799566837736	3
Chhattisgarh	Tea (black)	298.05765765765767	1
Chhattisgarh	Oil (groundnut)	179.65989583333328	2
Chhattisgarh	Oil (mustard)	165.59371681415925	3
Delhi	Tea (black)	221.97128205128203	1
Delhi	Oil (groundnut)	171.78246376811592	2
Delhi	Oil (sunflower)	129.1298959072165	3
Goa	Tea (black)	241.61477477477482	1
Goa	Oil (groundnut)	178.1410112359596	2
Goa	Oil (mustard)	158.98766233766233	3
Gujarat	Tea (black)	265.50549242424233	1
Gujarat	Oil (groundnut)	148.2668120808569	2
Gujarat	Oil (soybean)	126.1780386158615865	3
Haryana	Tea (black)	213.17700000000000	1
Haryana	Oil (groundnut)	159.06504505926304	2
Haryana	Oil (sunflower)	156.97300000000000	3
Himachal Pradesh	Tea (black)	233.442397638903686	1
Himachal Pradesh	Oil (groundnut)	167.89349286349296	2
Himachal Pradesh	Oil (sunflower)	166.56114942528734	3
Jharkhand	Tea (black)	258.7871578947369	1
Jharkhand	Oil (groundnut)	177.88599999999999	2
Jharkhand	Oil (sunflower)	157.0774537837837	3
Karnataka	Tea (black)	286.65858808778586	1
Karnataka	Oil (mustard)	164.5722434367541	2
Karnataka	Oil (groundnut)	156.9955649717515	3
Kerala	Tea (black)	294.70535087719293	1
Kerala	Oil (groundnut)	166.3139837398374	2
Kerala	Oil (mustard)	148.5885893999625	3
Madhya Pradesh	Tea (black)	247.42497326283215	1
Madhya Pradesh	Oil (groundnut)	158.56956417910442	2
Madhya Pradesh	Oil (sunflower)	128.1180962800875	3
Maharashtra	Tea (black)	242.90563157894738	1
Maharashtra	Oil (groundnut)	154.1169819607842	2
Maharashtra	Oil (sunflower)	121.15905707196026	3
Manipur	Tea (black)	263.995	1
Manipur	Oil (soybean)	179.18789473684208	2
Manipur	Oil (sunflower)	178.14888888888888	3
Meghalaya	Tea (black)	287.58265193370164	1
Meghalaya	Oil (sunflower)	148.70181818181814	2
Meghalaya	Oil (soybean)	128.4872641589434	3
Mizoram	Tea (black)	268.6525454545455	1
Mizoram	Oil (groundnut)	171.77599999999998	2
Mizoram	Ghee (vanaspati)	168.62	3

## Relevance

Ranking commodities by price within each state creates a clear picture of local food affordability. Identifying the top 3 most expensive items per state reveals which essential foods might be out of reach for lower-income households in those regions.

Q9) How do food prices in metropolitan cities compare to other markets during different seasons of the year?

```

SELECT
CASE
    WHEN city IN ('Delhi', 'Mumbai', 'Kolkata', 'Chennai', 'Bangalore', 'Hyderabad')
THEN 'Metropolitan'
    ELSE 'Other'
END AS city_type,
season,
category,
AVG(price) AS avg_price
FROM food_prices
GROUP BY
CASE
    WHEN city IN ('Delhi', 'Mumbai', 'Kolkata', 'Chennai', 'Bangalore', 'Hyderabad')
THEN 'Metropolitan'
    ELSE 'Other'
END,
season,
category
ORDER BY
category,
season,
city_type;

```

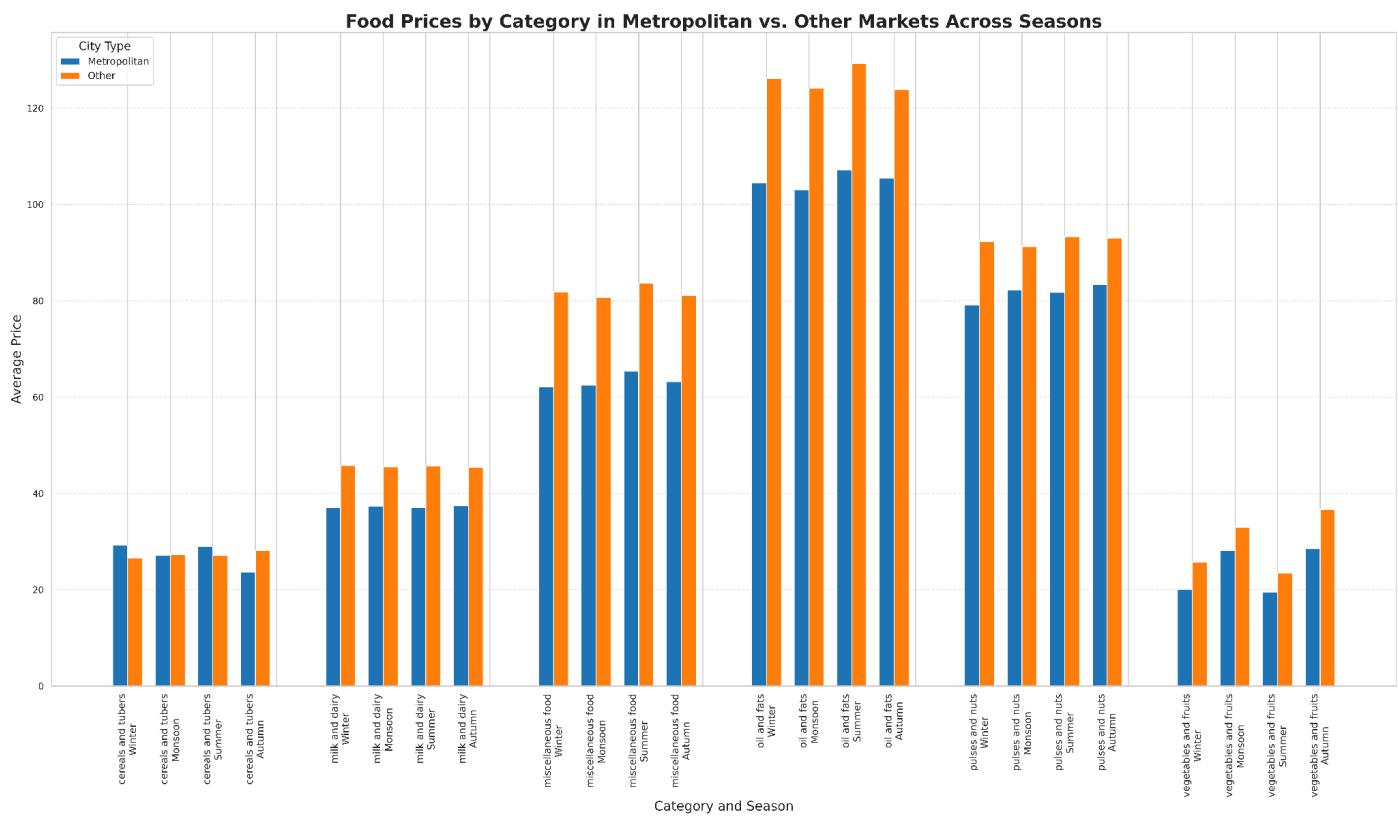
```

0: jdbc:hive2://localhost:10000>
0: jdbc:hive2://localhost:10000>
0: jdbc:hive2://localhost:10000> SELECT CASE WHEN city IN ('Delhi', 'Mumbai', 'Kolkata', 'Chennai', 'Bangalore', 'Hyderabad') THEN 'Metropolitan' ELSE 'Other' END AS city_type, season, category, AVG(price)
AS avg_price FROM food_prices GROUP BY CASE WHEN city IN ('Delhi', 'Mumbai', 'Kolkata', 'Chennai', 'Bangalore', 'Hyderabad') THEN 'Metropolitan' ELSE 'Other' END, season, category ORDER BY category, sea
son, city_type;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
+-----+-----+-----+-----+
| city_type | season | category | avg_price |
+-----+-----+-----+-----+
| Metropolitan | Autumn | cereals and tubers | 23.663690036900366 |
| Other | Autumn | cereals and tubers | 28.16463473744943 |
| Metropolitan | Monsoon | cereals and tubers | 27.077121019108255 |
| Other | Monsoon | cereals and tubers | 27.334542167269294 |
| Metropolitan | Summer | cereals and tubers | 28.940518075376887 |
| Other | Summer | cereals and tubers | 27.08474465691793 |
| Metropolitan | Winter | cereals and tubers | 29.21227784738911 |
| Other | Winter | cereals and tubers | 26.5883796363754 |
| Metropolitan | Autumn | milk and dairy | 37.483172431939311 |
| Other | Autumn | milk and dairy | 45.4677858957795 |
| Metropolitan | Monsoon | milk and dairy | 37.3030041243457 |
| Other | Monsoon | milk and dairy | 45.445687448239234 |
| Metropolitan | Summer | milk and dairy | 37.0299647482014 |
| Other | Summer | milk and dairy | 45.66130142540643 |
| Metropolitan | Winter | milk and dairy | 36.99626865671442 |
| Other | Winter | milk and dairy | 45.74775538392286 |
| Metropolitan | Autumn | miscellaneous food | 63.17900567851779 |
| Other | Autumn | miscellaneous food | 81.0918886476823 |
| Metropolitan | Monsoon | miscellaneous food | 62.4566851851852 |
| Other | Monsoon | miscellaneous food | 88.66133437785601 |
| Metropolitan | Summer | miscellaneous food | 65.3922014925373 |
| Other | Summer | miscellaneous food | 83.63313759641716 |
| Metropolitan | Winter | miscellaneous food | 62.124851301115235 |
| Other | Winter | miscellaneous food | 81.77023647765891 |
| Metropolitan | Autumn | oil and fats | 105.46684729864043 |
| Other | Autumn | oil and fats | 123.82246940153184 |
| Metropolitan | Monsoon | oil and fats | 102.996567164179 |
| Other | Monsoon | oil and fats | 124.18728844855544 |
| Metropolitan | Summer | oil and fats | 107.17266284881636 |
| Other | Summer | oil and fats | 129.27451437604952 |
| Metropolitan | Winter | oil and fats | 104.44171428571434 |
| Other | Winter | oil and fats | 126.11663161182857 |
| Metropolitan | Autumn | pulses and nuts | 83.35692307692308 |
| Other | Autumn | pulses and nuts | 92.99366506313889 |
| Metropolitan | Monsoon | pulses and nuts | 82.249904627962954 |
| Other | Monsoon | pulses and nuts | 91.22812328266362 |
| Metropolitan | Summer | pulses and nuts | 81.72578199052133 |
| Other | Summer | pulses and nuts | 93.28789627770843 |
| Metropolitan | Winter | pulses and nuts | 79.11886513994988 |
| Other | Winter | pulses and nuts | 92.23735688936891 |
| Metropolitan | Autumn | vegetables and fruits | 28.4667547892723 |
| Other | Autumn | vegetables and fruits | 34.487844677651765 |
| Metropolitan | Monsoon | vegetables and fruits | 28.100165642913906 |
| Other | Monsoon | vegetables and fruits | 32.9844231292517 |
| Metropolitan | Summer | vegetables and fruits | 19.4812772727272735 |
| Other | Summer | vegetables and fruits | 23.451075067824156 |
| Metropolitan | Winter | vegetables and fruits | 28.007400759124093 |
| Other | Winter | vegetables and fruits | 25.663752380952396 |
+-----+
8 rows selected (4.879 seconds)
0: jdbc:hive2://localhost:10000>
0: jdbc:hive2://localhost:10000>
0: jdbc:hive2://localhost:10000>
0: jdbc:hive2://localhost:10000>
```

## Relevance

Comparing metropolitan and non-metropolitan food prices across seasons reveals urban-rural divides and seasonal vulnerabilities. This analysis helps understand whether urban consumers consistently pay premiums regardless of season, or if there are specific times when the urban-rural price gap widens or narrows, which impacts migration patterns and urban poverty levels.

## Visualization



Q10) Analyse the evolution of commodity prices over time, revealing year-to-year price dynamics and identifying significant market trends.

```

SELECT
    commodity,
    year,
    AVG(price) AS avg_price,
    LAG(AVG(price), 1) OVER (PARTITION BY commodity ORDER BY year) AS
    prev_year_price,
    ((AVG(price) - LAG(AVG(price), 1) OVER (PARTITION BY commodity ORDER BY year)) /
    LAG(AVG(price), 1) OVER (PARTITION BY commodity ORDER BY year)) * 100 AS pct_change
FROM
    food_prices
GROUP BY
    commodity, year
ORDER BY
    commodity, year;

```

commodity	year	avg_price	prev_year_price	pct_change
Chickpeas	2000	20.351785714285715	NULL	NULL
Chickpeas	2001	25.371153846153845	20.351785714285715	24.663035481671958
Chickpeas	2002	23.17241379310345	25.371153846153845	-8.66629980391273
Chickpeas	2003	22.31864846486467	23.17241379310345	-3.6844031531531596
Chickpeas	2004	22.889189189189189	22.31864846486467	-0.13326577878887634
Chickpeas	2005	24.29259259259259	22.88918918918917	8.989550731296116
Chickpeas	2006	32.57441764705886	24.29259259259259	34.0195262887748
Chickpeas	2007	35.98216216216216	32.57441764705886	4.33498394964574
Chickpeas	2008	35.53484848484848	33.98216216216214	4.55664807857031
Chickpeas	2009	34.349800000000004	35.53484848484849	-3.37142369848075
Chickpeas	2010	34.37694444444444	34.349800000000004	0.08135446285826217
Chickpeas	2011	43.132988645161595	34.37694444444444	25.46958975767605
Chickpeas	2012	65.12325744444444	43.132988645161595	41.46958975767605
Chickpeas	2013	55.36124999999999	61.17241379310345	-9.53234214459981
Chickpeas	2014	109.296388888888885	55.36124999999999	14.92288502899993
Chickpeas	2015	61.40432432432433	49.296388888888885	24.515058391515184
Chickpeas	2016	182.86764705882352	61.40432432432433	66.2224649282847
Chickpeas	2017	182.86764705882352	182.86764705882352	-15.445638705557413
Chickpeas	2018	66.83898823529413	86.38264705882352	-22.562527982831672
Chickpeas	2019	66.66894736842104	66.83898823529413	4.247095848566667
Chickpeas	2020	73.71827692387692	69.66894736842184	8.511957418949397
Ghee (vanaspati)	2015	73.74183266932273	NULL	NULL
Ghee (vanaspati)	2016	75.17497164461251	73.74183266932273	1.94346515134189185
Ghee (vanaspati)	2017	78.717572815534	75.17497164461251	4.712474236330983
Ghee (vanaspati)	2018	80.92685042016895	78.717572815534	2.8055712656447397
Ghee (vanaspati)	2019	81.54829341317367	80.92685042016895	0.7689832820899674
Ghee (vanaspati)	2020	94.35707717569792	81.54829341317367	15.706991883773382
Ghee (vanaspati)	2021	129.52185761847468	94.35707717569792	37.267772049888865
Ghee (vanaspati)	2022	158.76835451977482	129.52185761847468	16.48972611743084
Ghee (vanaspati)	2023	133.78655172413784	150.76835451977482	-11.263506691663988
Ghee (vanaspati)	2024	127.25	133.78655172413784	-4.885806263708723
Lentils	2013	51.81151898734176	NULL	NULL
Lentils	2014	47.10025125628143	51.81151898734176	-9.09380838078298
Lentils	2015	58.396993736952	47.10025125628143	23.984463308828912
Lentils	2016	92.01198294243065	58.396993736952	57.56287619341614
Lentils	2017	85.8033264033264	92.01198294243065	-6.74660837815916
Lentils	2018	64.62271428571425	85.8033264033264	-24.681071086938517
Lentils	2019	81.5868728522336	64.62271428571425	26.24158538867956
Lentils	2020	97.82982935153578	81.5868728522336	19.91787836938323
Lentils	2021	105.51795180722891	97.82982935153578	7.8586689833293875
Lentils	2022	106.3027891156458	105.51795180722891	0.882693665481886
Lentils	2023	108.889394495412827	106.3027891156458	1.29362288951483
Lentils	2024	149.992	108.889394495412827	24.77613851767933
Lentils (masur)	2012	52.640496183206995	NULL	NULL
Lentils (masur)	2013	57.222027397240556	52.640496183206995	8.703434705685403
Lentils (masur)	2014	66.7657808788485	57.222027397240556	16.67844924449238
Lentils (masur)	2015	82.39393406593499	66.7657808788485	23.49743036919912
Lentils (masur)	2016	82.5458174976644	82.39393406593499	0.184380165714415
Lentils (masur)	2017	68.66752866115783	82.5458174976644	-16.812849655167314
Lentils (masur)	2018	68.76174568965517	68.66752866115783	-11.51312133567941
Lentils (masur)	2019	63.17878942536823	69.76174568965517	3.977946438587581
Lentils (masur)	2020	75.57837795275589	63.17878942536823	19.62618847388999
Lentils (masur)	2021	88.16298856155587	75.57837795275589	16.65186204934146
Lentils (masur)	2022	95.997392996119893	88.16298856155587	8.88623545148517
Lentils (masur)	2023	91.77799783549781	95.997392996119893	-4.39532652962289
Lentils (masur)	2024	93.795	91.77799783549781	2.197696846816649
Lentils (moong)	2015	102.10896551724146	NULL	NULL

## Relevance

Tracking year-to-year price dynamics exposes longer-term food security threats or improvements. By analyzing how prices evolve over time, we can separate random fluctuations from genuine trends, helping forecast future prices and identify which commodities are becoming increasingly unaffordable or more accessible over time.

Q11) Identify the top 5 markets with the highest average price for each commodity, considering only the latest year.

```
WITH latest_year AS (
    SELECT MAX(year) AS max_year FROM food_prices
),
ranked_data AS (
    SELECT
        fp.year, fp.state, fp.market, fp.commodity,
        AVG(fp.price) AS avg_price,
        RANK() OVER (PARTITION BY fp.commodity ORDER BY AVG(fp.price) DESC) AS market_rank
    FROM food_prices fp
    JOIN latest_year ly ON fp.year = ly.max_year
    GROUP BY fp.year, fp.state, fp.market, fp.commodity
)
SELECT
    year, state, market, commodity, avg_price, market_rank
FROM ranked_data WHERE market_rank <= 5
ORDER BY commodity, market_rank, state, year;
```

year	state	market	commodity	avg_price	market_rank
2024	Assam	North-East Zone	Ghee (vanaspati)	143.88	1
2024	Rajasthan	North Zone	Ghee (vanaspati)	127.82	2
2024	Karnataka	South Zone	Ghee (vanaspati)	125.15	3
2024	Maharashtra	West Zone	Ghee (vanaspati)	122.68	4
2024	Orissa	East Zone	Ghee (vanaspati)	116.72	5
2024	Karnataka	South Zone	Lentils	145.57	1
2024	Maharashtra	West Zone	Lentils	148.04	2
2024	Orissa	East Zone	Lentils	147.64	3
2024	Rajasthan	North Zone	Lentils	147.53	4
2024	Assam	North-East Zone	Lentils	141.18	5
2024	Karnataka	South Zone	Lentils (masur)	101.14	1
2024	Assam	North-East Zone	Lentils (masur)	95.89	2
2024	Maharashtra	West Zone	Lentils (masur)	91.54	3
2024	Orissa	East Zone	Lentils (masur)	86.61	4
2024	Karnataka	South Zone	Lentils (moong)	124.26	1
2024	Assam	North-East Zone	Lentils (moong)	118.29	2
2024	Orissa	East Zone	Lentils (moong)	115.17	3
2024	Maharashtra	West Zone	Lentils (moong)	112.64	4
2024	Rajasthan	North Zone	Lentils (moong)	112.64	4
2024	Karnataka	South Zone	Lentils (urad)	140.15	1
2024	Orissa	East Zone	Lentils (urad)	118.47	2
2024	Maharashtra	West Zone	Lentils (urad)	116.7	3
2024	Assam	North-East Zone	Milk (pasteurized)	71.83	1
2024	Rajasthan	North Zone	Milk (pasteurized)	69.3	2
2024	Maharashtra	West Zone	Milk (pasteurized)	57.1	3
2024	Orissa	East Zone	Milk (pasteurized)	69.15	4
2024	Karnataka	South Zone	Milk (pasteurized)	69.0	5
2024	Rajasthan	North Zone	Oil (groundnut)	192.87	1
2024	Karnataka	South Zone	Oil (groundnut)	192.75	2
2024	Assam	North-East Zone	Oil (groundnut)	192.01	3
2024	Maharashtra	West Zone	Oil (groundnut)	189.81	4
2024	Orissa	East Zone	Oil (groundnut)	185.04	5
2024	Assam	North-East Zone	Oil (mustard)	149.64	1
2024	Orissa	East Zone	Oil (mustard)	133.0	2
2024	Assam	North-East Zone	Oil (palm)	118.25	1
2024	Rajasthan	North Zone	Oil (palm)	109.58	2
2024	Maharashtra	West Zone	Oil (palm)	105.68	3
2024	Orissa	East Zone	Oil (palm)	99.95	4
2024	Karnataka	South Zone	Oil (palm)	93.78	5
2024	Karnataka	South Zone	Oil (soybean)	137.5	1
2024	Assam	North-East Zone	Oil (soybean)	136.11	2

## Relevance

Identifying markets with the highest prices for each commodity in the latest year reveals current affordability hotspots. This time-sensitive information shows where consumers are currently facing the greatest pressure, helping target immediate relief efforts and understand whether high-price areas are persistent or shifting over time.

Q12) Analyze the impact of seasonality on price volatility to calculate the coefficient of variation for each commodity.

```

SELECT
    commodity,
    season,
    AVG(price) AS avg_price,
    STDDEV(price) AS stddev_price,
    (STDDEV(price) / AVG(price)) * 100 AS coeff_variation
FROM
    food_prices
GROUP BY
    commodity, season
ORDER BY
    commodity, season;

```

SELECT commodity, season, AVG(price) AS avg_price, STDDEV(price) AS stddev_price, (STDDEV(price) / AVG(price)) * 100 AS coeff_variation FROM food_prices GROUP BY commodity, season ORDER BY commodity, season;				
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.				
commodity	season	avg_price	stddev_price	coeff_variation
Chickpeas	Autumn	49.75468085106375	28.042488075829223	56.361748311451706
Chickpeas	Monsoon	46.909523163742696	24.290053021854867	51.7863370375953
Chickpeas	Summer	44.472987804878066	21.59098318496481	48.54853287504237
Chickpeas	Winter	46.80564780582355	21.80289842831527	46.5817697780854
Ghee (vanaspati)	Autumn	113.95461485679807	33.5624008777614	29.45242798784849
Ghee (vanaspati)	Monsoon	112.47162805662812	36.790043472395915	32.710518293262624
Ghee (vanaspati)	Summer	116.17740543735212	37.04074238547987	31.882483729121984
Ghee (vanaspati)	Winter	113.03012820512828	33.152888511413764	29.331819116644324
Lentils	Autumn	94.2528865241845	25.48180666802516	27.835588228862693
Lentils	Monsoon	91.7355108645653	26.11772144431787	28.478677339987
Lentils	Summer	92.26643280900805	26.1302794313932	28.32845158686584
Lentils	Winter	92.6594614939899	24.66129542916659	26.912858564291723
Lentils (masur)	Autumn	81.350369517051702	17.02562578406364	21.38021666806993
Lentils (masur)	Monsoon	80.350369517051702	17.02119803749384	21.37976777152134
Lentils (masur)	Summer	80.4702481751827	17.318351176742798	21.46307722619976
Lentils (masur)	Winter	80.50281035039547	17.05440769881658	21.184860926421244
Lentils (moong)	Autumn	94.67383984676	15.765002929921772	16.3974136142721248
Lentils (moong)	Monsoon	95.2164227642277	16.36387753867288	17.185124283443636
Lentils (moong)	Summer	99.73579008585242	15.755389380069254	15.797141134834722
Lentils (moong)	Winter	97.42632936507931	14.967988232995317	15.36523873877725
Lentils (urad)	Autumn	104.39674498426195	21.4425494252237168	20.539488681989987
Lentils (urad)	Monsoon	101.62187443902442	22.43585797278577	22.87796871504443
Lentils (urad)	Summer	104.41640468227428	21.4967588996178052	20.58752948024779
Lentils (urad)	Winter	103.64351145838157	19.8254997899403	19.12578288724685
Milk	Autumn	35.532849999999915	8.138834438673918	22.9851774793968
Milk	Monsoon	35.38723958333333	8.53121957777827	24.10818045778359
Milk	Summer	35.00424581085587	8.684899298658884	24.8109988203505083
Milk	Winter	33.722040816326526	8.717280452533918	25.85039410874993
Milk (pasteurized)	Autumn	45.650866225165546	9.719827247385814	21.2928331798936
Milk (pasteurized)	Monsoon	45.77698764169661	9.79851298569446	21.39899951593173
Milk (pasteurized)	Summer	46.00653658536583	10.529777465273124	22.887568268512204
Milk (pasteurized)	Winter	46.365887265135676	9.885561176113857	21.32863411999373
Oil (groundnut)	Autumn	153.94389819156058	33.56972864832424	21.886469137576108
Oil (groundnut)	Monsoon	153.7168546195653	33.98183952736364	22.18677520484963
Oil (groundnut)	Summer	160.67317934782636	34.8125618912947	21.66669137475122
Oil (groundnut)	Winter	156.364309961547474	33.17056533763519	21.213642382114986
Oil (mustard)	Autumn	119.11684647382983	48.22867922773225	40.48854694176614
Oil (mustard)	Monsoon	117.6244642188999	47.35378186078775	40.258446383215286
Oil (mustard)	Summer	121.9378987778587	49.360276223762014	40.47997748773678
Oil (mustard)	Winter	121.98881843464443	49.21936829225634	40.3499186158158
Oil (palm)	Autumn	108.36812461828418	28.14379964768855	28.848575385696384
Oil (palm)	Monsoon	102.12359865298491	30.43967816428587	30.7859842898594
Oil (palm)	Summer	102.9246887930810	32.000000000000004	30.88487816428587
Oil (palm)	Winter	102.9246887930810	28.493524458911433	27.6838571119341
Oil (soybean)	Autumn	121.6812977899247	33.51874573617435	27.532898857773733
Oil (soybean)	Monsoon	123.6226608187152	35.53289885307518	28.742468727098132
Oil (soybean)	Summer	127.45671184922824	36.446667525584808	28.5951492926708196
Oil (soybean)	Winter	124.34885074426871	31.621439475899873	25.42961936288013
Oil (sunflower)	Autumn	138.58737517592484	37.1967381912985	28.484789235495662
Oil (sunflower)	Monsoon	132.0285946882126	48.24712353989511	30.48548874285649
Oil (sunflower)	Summer	138.36912429378535	48.78388418521534	29.4746599932785573
Oil (sunflower)	Winter	133.11793522267217	35.344641650794554	26.551374697541824
Onions	Autumn	35.2285943889616	17.56016544998633	49.85763547807356
Onions	Monsoon	26.1885430463576	10.9897289898907981	41.963888509929484
Onions	Summer	28.67624466457738	7.828883855729173	37.82545336233451
Onions	Winter	27.431321671525744	11.943148881996658	43.53836444786886
Potatoes	Autumn	24.414113345521	10.34881081905663	42.388641051362406

## Relevance

**High CV:** A high coefficient of variation indicates that a commodity's price is highly volatile relative to its average price. This suggests that the commodity is more susceptible to seasonal changes or external factors affecting its price.

**Low CV:** A low CV suggests that the commodity's price is relatively stable, with less fluctuation around the mean. This can indicate a more predictable pricing pattern, which is beneficial for planning and forecasting.

Q13) Identify the cheapest and costliest state for each commodity for a given year? Take the year 2024 for example.

```
WITH price_extremes AS (
    SELECT
        commodity,
        state,
        MIN(price) AS min_price,
        MAX(price) AS max_price
    FROM food_prices
    WHERE year = 2024
    GROUP BY commodity, state
),
ranked_extremes AS (
    SELECT
        commodity,
        state AS cheapest_state,
        min_price AS cheapest_price,
        NULL AS costliest_state,
        NULL AS costliest_price,
        RANK() OVER (PARTITION BY commodity ORDER BY min_price ASC) AS cheapest_rank,
        RANK() OVER (PARTITION BY commodity ORDER BY max_price DESC) AS
costliest_rank
    FROM price_extremes
    UNION ALL
    SELECT
        commodity,
        NULL AS cheapest_state,
        NULL AS cheapest_price,
        state AS costliest_state,
        max_price AS costliest_price,
        RANK() OVER (PARTITION BY commodity ORDER BY min_price ASC) AS cheapest_rank,
        RANK() OVER (PARTITION BY commodity ORDER BY max_price DESC) AS
costliest_rank
    FROM price_extremes
)
SELECT
    commodity,
    MAX(CASE WHEN cheapest_rank = 1 THEN cheapest_state END) AS cheapest_state,
    MAX(CASE WHEN cheapest_rank = 1 THEN cheapest_price END) AS cheapest_state_price,
    MAX(CASE WHEN costliest_rank = 1 THEN costliest_state END) AS costliest_state,
    MAX(CASE WHEN costliest_rank = 1 THEN costliest_price END) AS
costliest_state_price
FROM ranked_extremes
GROUP BY commodity
ORDER BY commodity;
```

```

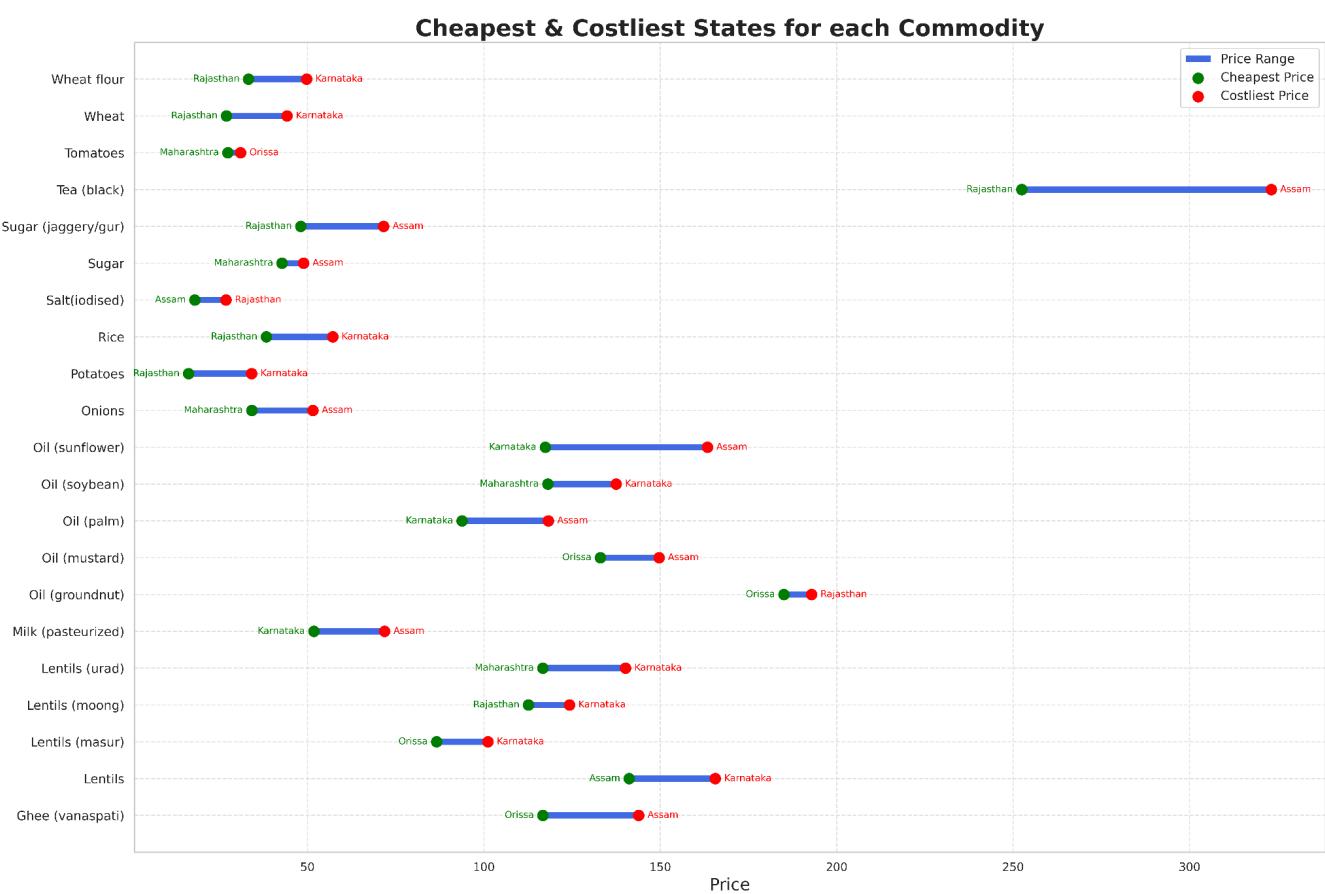
0: jdbc:hive2://localhost:10000> WITH price_extremes AS ( SELECT commodity, state, MIN(price) AS min_price, MAX(price) AS max_price FROM food_prices WHERE year = 2024 GROUP BY commodity, state ), ranked_extremes AS (SELECT commodity, state AS cheapest_state, min_price AS cheapest_price, NULL AS costliest_state, NULL AS costliest_price, RANK() OVER (PARTITION BY commodity ORDER BY min_price ASC) AS cheapest_rank, RANK() OVER (PARTITION BY commodity ORDER BY max_price DESC) AS costliest_rank FROM price_extremes UNION ALL SELECT commodity, NULL AS cheapest_state, NULL AS cheapest_price, state AS costliest_state, max_price AS costliest_price, RANK() OVER (PARTITION BY commodity ORDER BY min_price ASC) AS cheapest_rank, RANK() OVER (PARTITION BY commodity ORDER BY max_price DESC) AS costliest_rank FROM price_extremes ) SELECT commodity, MAX(CASE WHEN cheapest_rank = 1 THEN cheapest_state END) AS cheapest_state, MAX(CASE WHEN cheapest_rank = 1 THEN cheapest_price END) AS cheapest_state_price, MAX(CASE WHEN costliest_rank = 1 THEN costliest_state END) AS costliest_state, MAX(CASE WHEN costliest_rank = 1 THEN costliest_price END) AS costliest_state_price FROM ranked_extremes GROUP BY commodity ORDER BY commodity
;
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
+-----+-----+-----+-----+-----+
| commodity | cheapest_state | cheapest_state_price | costliest_state | costliest_state_price |
+-----+-----+-----+-----+-----+
| Ghee (vanaspati) | Orissa | 116.72 | Assam | 143.88 |
| Lentils | Assam | 141.18 | Karnataka | 165.57 |
| Lentils (masur) | Orissa | 86.61 | Karnataka | 181.14 |
| Lentils (moong) | Maharashtra | 112.64 | Karnataka | 124.26 |
| Lentils (urad) | Maharashtra | 116.2 | Karnataka | 141.15 |
| Milk (pasteurized) | Karnataka | 51.8 | Assam | 71.83 |
| Oil (groundnut) | Orissa | 105.04 | Rajasthan | 122.87 |
| Oil (mustard) | Orissa | 133.0 | Assam | 149.44 |
| Oil (palm) | Karnataka | 93.78 | Assam | 118.25 |
| Oil (soybean) | Maharashtra | 118.11 | Karnataka | 137.5 |
| Oil (sunflower) | Karnataka | 117.42 | Assam | 163.38 |
| Onions | Maharashtra | 34.21 | Assam | 51.54 |
| Potatoes | Rajasthan | 16.27 | Karnataka | 34.16 |
| Rice | Rajasthan | 38.32 | Karnataka | 57.15 |
| Salt (iodised) | Assam | 18.05 | Rajasthan | 26.89 |
| Sugar | Maharashtra | 42.78 | Assam | 48.9 |
| Sugar (jaggery/gur) | Rajasthan | 48.1 | Assam | 71.58 |
| Tea (black) | Rajasthan | 252.46 | Assam | 323.2 |
| Tomatoes | Maharashtra | 27.41 | Orissa | 31.03 |
| Wheat | Rajasthan | 27.04 | Karnataka | 44.19 |
| Wheat flour | Rajasthan | 33.28 | Karnataka | 49.74 |
+-----+-----+-----+-----+-----+
21 rows selected (13.208 seconds)
0: jdbc:hive2://localhost:10000>

```

## Relevance

Comparing the cheapest and costliest states for each commodity reveals the extremes of India's regional price disparities. This stark contrast highlights opportunities for improving inter-state trade and distribution networks, while also helping consumers understand which foods might be worth purchasing from neighboring states when possible.

## Visualization



Q14) Which cities experienced consistent food price inflation over the years?

```
SELECT
    city,
    COUNT(DISTINCT year) AS years_with_inflation
FROM (
    SELECT
        city,
        year,
        AVG(price) AS avg_yearly_price,
        LAG(AVG(price)) OVER (PARTITION BY city ORDER BY year) AS prev_year_price
    FROM food_prices
    GROUP BY city, year
) t
WHERE t.avg_yearly_price > t.prev_year_price
GROUP BY city
ORDER BY years_with_inflation DESC;
```

city	years_with_inflation
Lucknow	22
Kamrup	22
Shimla	21
Khordha	21
Bikaner	21
Bangalore Urban	21
Chennai	20
Hyderabad	20
Mumbai City	20
Thiruvananthapuram	20
East Khasi Hills	19
Ahmedabad	19
Delhi	19
Dibrugarh	19
Pune	18
Aizawl	18
West Tripura	18
Bhopal	17
Amravati	11
Ludhiana	11
Amritsar	10
Sambalpur	10
Ernakulam	10
Ranchi	10
Hisar	10
Bathinda	10
Chandigarh	10
Dehra Dun	10
Dharwad	10
Kota	9
Krishna	9
Banksa	9
Bid	9
Agra	9
Cuttack	9
Dimapur	9
Dindigul	9
Kanpur	9
Karnal	9
Tiruchirappalli	8
Indore	8
Darjiling	8
Varanasi	8
Puducherry	8
Majkot	8
Mandi	7
Gopalpur	7
Andaman Islands	7
Gurgaon	6
North Goa	6
Gwalior	6
Kozhikode	5
Jabalpur	5
Panchkula	5
Sundargarh	5
Kohima	3
Udaipur	3

## Relevance

Tracking cities with consistent food price inflation identifies urban areas facing chronic affordability challenges. These locations may have structural issues that continuously drive prices upward, such as rapid population growth, limited agricultural land in surrounding areas, or inefficient distribution systems that need targeted policy interventions.

Q15) Identify the top 10 regions in India where food prices are the most volatile, by calculating the average price variance across nearby markets (within a 100 km radius).

```

WITH markets AS (
    SELECT
        state, city, market,
        AVG(price) AS avg_price,
        MAX(latitude) AS latitude,
        MAX(longitude) AS longitude,
        ROW_NUMBER() OVER (PARTITION BY state, city, market) AS rn
    FROM food_prices
    GROUP BY state, city, market
),
nearby_markets AS (
    SELECT
        m1.state, m1.city, m1.market, m1.avg_price AS price1, m2.market AS nearby_market,
        m2.avg_price AS price2,
        111.045 * DEGREES(ACOS(
            LEAST(1.0, COS(RADIANS(m1.latitude)) * COS(RADIANS(m2.latitude)) *
            COS(RADIANS(m2.longitude) - RADIANS(m1.longitude)) +
            SIN(RADIANS(m1.latitude)) * SIN(RADIANS(m2.latitude))))
        )) AS distance_km
    FROM markets m1 JOIN markets m2 ON m1.market != m2.market
    WHERE
        111.045 * DEGREES(ACOS(
            LEAST(1.0, COS(RADIANS(m1.latitude)) * COS(RADIANS(m2.latitude)) *
            COS(RADIANS(m2.longitude) - RADIANS(m1.longitude)) +
            SIN(RADIANS(m1.latitude)) * SIN(RADIANS(m2.latitude))))
        )) <= 100
)
SELECT
    state, city, market,
    AVG(ABS(price1 - price2)) AS avg_price_volatility,
    COUNT(DISTINCT nearby_market) AS nearby_market_count
FROM nearby_markets
GROUP BY state, city, market
HAVING COUNT(DISTINCT nearby_market) > 2
ORDER BY avg_price_volatility DESC LIMIT 10;

```

```

0: jdbc:hive2://localhost:10000> WITH markets AS ( SELECT state, city, market, AVG(price) AS avg_price, MAX(latitude) AS latitude, MAX(longitude) AS longitude, ROW_NUMBER() OVER (PARTITION BY state, city, market) AS rn FROM food_prices GROUP BY state, city, market ), nearby_markets AS ( SELECT m1.state, m1.city, m1.market, m1.avg_price AS price1, m2.market AS nearby_market, m2.avg_price AS price2, 111.045 * DEGREES(ACOS(LEAST(1.0, COS(RADIANS(m1.latitude)) * COS(RADIANS(m2.latitude)) * COS(RADIANS(m2.longitude) - RADIANS(m1.longitude)) + SIN(RADIANS(m1.latitude)) * SIN(RADIANS(m2.latitude)))))) AS distance_km FROM markets m1 JOIN markets m2 ON m1.market != m2.market WHERE 111.045 * DEGREES(ACOS(LEAST(1.0, COS(RADIANS(m1.latitude)) * COS(RADIANS(m2.latitude)) * COS(RADIANS(m2.longitude) - RADIANS(m1.longitude)) + SIN(RADIANS(m1.latitude)) * SIN(RADIANS(m2.latitude)))))) <= 100 ) SELECT state, city, market, AVG(ABS(price1 - price2)) AS avg_price_volatility, COUNT(DISTINCT nearby_market) AS nearby_market_count FROM nearby_markets GROUP BY state, city, market HAVING COUNT(DISTINCT nearby_market) > 2 ORDER BY avg_price_volatility DESC LIMIT 10;
WARNING: Hive-on-MR is deprecated and will be removed in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/opt/hive/lib/lon4j-slf4j-impl-2.6.2.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/opt/hadoop-2.7.1/share/hadoop/common/lib/slf4j-log4j12-1.7.10.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Execution log at: /tmp/root/root_202503231046:40_a3427828-823c-4942-b25c-e34678988187.log
2025-03-23 10:46:40 Starting to launch local task to process map join; maximum memory = 477626368
2025-03-23 10:46:40 Dump the side-table for tag: 1 with group count: 1 into file: file:/tmp/root/135576f2-db96-4567-bd79-7b7c57dfcc70/hive_2025-03-23_10-46-30_811_5025985908235006507-1-local-10008/Ha
shTable-Stage-6/MapJoin-mapfile121--.hashtable
2025-03-23 10:46:41 Uploaded 1 File to: file:/tmp/root/135576f2-db96-4567-bd79-7b7c57dfcc70/hive_2025-03-23_10-46-30_811_5025985908235006507-1-local-10008/HashTable-Stage-6/MapJoin-mapfile121--.hasht
able (6736 bytes)
2025-03-23 10:46:41 End of local task; Time Taken: 0.224 sec.
+-----+-----+-----+-----+-----+
| state | city | market | avg_price_volatility | nearby_market_count |
+-----+-----+-----+-----+-----+
| Karnataka | Bangalore Urban | Bengaluru | 31.7777886046837 | 4 |
| Assam | Kamrup | Guwahati | 29.682731774791034 | 3 |
| Meghalaya | East Khasi Hills | Shillong | 25.650885913298822 | 3 |
| Bihar | Banka | Bhagalpur | 23.12405612351834 | 4 |
| Jharkhand | Ranchi | Ranchi | 20.426011016500635 | 3 |
| Kerala | Thrissur | Thrissur | 20.0835545667978 | 4 |
| Himachal Pradesh | Solan | Solan | 19.225279582612383 | 5 |
| Uttar Pradesh | Saharanpur | Saharanpur | 19.210837360623998 | 3 |
| Uttar Pradesh | Agra | Agra | 19.119423780261184 | 3 |
| Jharkhand | Gumla | Gumla | 16.196884757246558 | 4 |
+-----+-----+-----+-----+-----+
10 rows selected (14.268 seconds)
0: jdbc:hive2://localhost:10000>

```

## Relevance

We wanted to approach this problem using geospatial analysis to uncover localized price inconsistencies that wouldn't be visible through standard market analysis. Hive doesn't natively support geospatial functions, but you can use the [Haversine formula](#) to calculate the distance between two latitude-longitude pairs.

According to the formula, distance between two geographical points is:

$$d = 2r \cdot \arcsin\left(\sqrt{\sin^2\left(\frac{\Phi_2 - \Phi_1}{2}\right) + \cos(\Phi_1) \cdot \cos(\Phi_2) \cdot \sin^2\left(\frac{\lambda_2 - \lambda_1}{2}\right)}\right)$$

Which in SQL converts to:

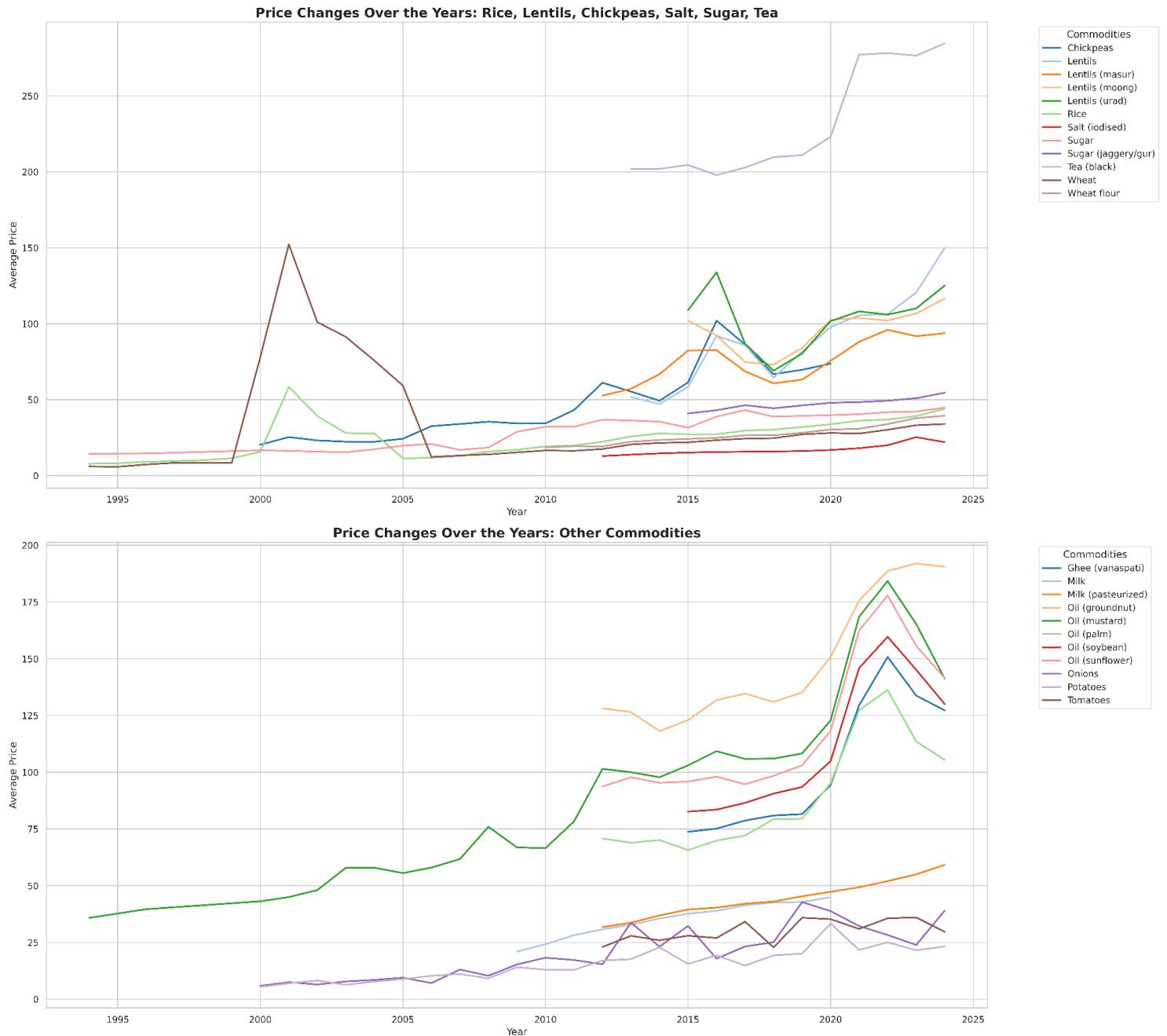
```
111.045 * DEGREES(ACOS(LEAST(1.0, COS(RADIANS(lat1)) * COS(RADIANS(lat2)) *  
    COS(RADIANS(long1) - RADIANS(long2)) +  
    SIN(RADIANS(lat1)) * SIN(RADIANS(lat2)))))
```

where 111.045 is a conversion factor representing approximately how many kilometers one degree of latitude equals at the equator. Reference - [Fast SQL location finder \(geolocation\) using Haversines](#)

Mapping regions with the highest food price volatility within a 100km radius reveals local market integration problems. Areas where nearby markets show wildly different prices likely suffer from poor transportation infrastructure, information gaps between markets, or fragmented supply chains. This spatial analysis helps prioritize logistics improvements that could stabilize regional food access.

# Conclusion

The analysis of Indian food prices reveals several key insights into market dynamics and regional variations. Monsoon emerges as the most expensive season. Notably, commodities such as tea, wheat, and ghee exhibit high volatility across different markets, with a *standard deviation* exceeding **45**, marking them as products with substantial price increases over the years. Tea and oil consistently rank as the most expensive products across nearly all states.



Bihar stands out with the most significant price difference between high and low-valued commodities, highlighting regional disparities in pricing. Interestingly, average prices across

seasons show that metropolitan areas tend to be cheaper compared to the rest of India, suggesting more competitive pricing in urban centers. Karnataka and Maharashtra consistently appear in the top five for the highest average price of each commodity, indicating their strong market presence.

Cities like Lucknow, Kamrup, Shimla, Khordha, Bikaner, Bengaluru, Chennai, and Mumbai City have experienced constant inflation over the past **22** years, reflecting ongoing economic pressures. Additionally, Bengaluru, Guwahati, and Shillong are particularly prone to price volatility, with significant price variance observed across nearby markets within a *100KM* range. These findings provide valuable insights for stakeholders aiming to navigate the complexities of the Indian food market.

## References

- [1] Configuring a MySQL Database for the Hive Metastore –  
<https://docs.ezmeral.hpe.com/datafabric-customer-managed/79/Hive/Config-RemoteMySQLForHiveMetastore.html>
- [2] Importing data to partitioned and clustered Hive tables –  
<https://www.ibm.com/docs/en/psfa/1.7.2?topic=movement-importing-data-partitioned-clustered-hive-tables>
- [3] Bubble Charts – <https://plotly.com/python/bubble-charts/>
- [4] Common Table Expression (CTE) –  
[https://hive.apache.org/docs/latest/common-table-expression\\_38572242/](https://hive.apache.org/docs/latest/common-table-expression_38572242/)
- [5] Hive Window Functions –  
<https://www.linkedin.com/pulse/hive-lag-rank-over-explained-example-fu-an-yang/>
- [6] Fast nearest-location finder for SQL –  
<https://www.plumislandmedia.net/mysql/haversine-mysql-nearest-loc/>