# Agricultural Crop Production & Yield Optimization Analytics System

End-to-End Data Engineering & Analytics Capstone

*Agneepradeep Verma*

# Automated Analytics Pipeline for Agricultural Intelligence

## Automated Pipeline

End-to-end agricultural analytics pipeline with minimal manual intervention

## Medallion Architecture

Bronze, Silver, Gold data layers with Airflow orchestration for production-grade reliability

## Interactive BI Dashboards

Power BI dashboards enabling real-time insights and data-driven decision making

**Technology Stack:** Python | Pandas | PySpark | Databricks | Apache Airflow | Power BI

# Business Context & Problem Statement

## Core Challenges

- **Data Volume & Distribution:** Agricultural data is large-scale, geographically distributed, and spans multiple dimensions (crops, regions, seasons, time)

- **Manual Processing Limitations:** Traditional manual analysis methods create bottlenecks and prevent timely, actionable insights

- **Scalability Requirements:** Growing data volumes demand automated, scalable analytics infrastructure

Solution Focus: Production-grade data pipeline enabling automated, enterprise-level agricultural analytics

# Medallion Architecture Implementation

## Bronze Layer: Raw Ingestion

- Python & Pandas for initial data loading
- Schema standardization across source files
- File-based storage with metadata tracking

## Silver Layer: Transformation & Cleansing

- Data cleaning and quality validation
- Year parsing and temporal alignment
- Yield recalculation and unit normalization
- Analysis-ready normalized dataset

## Gold Layer: Star Schema Modeling

- Fact and dimension table creation
- Star schema optimization for BI tools
- Performance-tuned aggregations

# Airflow Workflow Orchestration

## Docker-Based Orchestration

Apache Airflow provides robust pipeline management
with clear dependency handling and modular task design.

**01** **ingest_raw_to_bronze**

Raw data ingestion task

**02** **transform_bronze_to_silver**

Data cleaning and validation

**03** **transform_silver_to_gold**

Star schema generation

**Manual trigger DAG** enables controlled execution and
testing

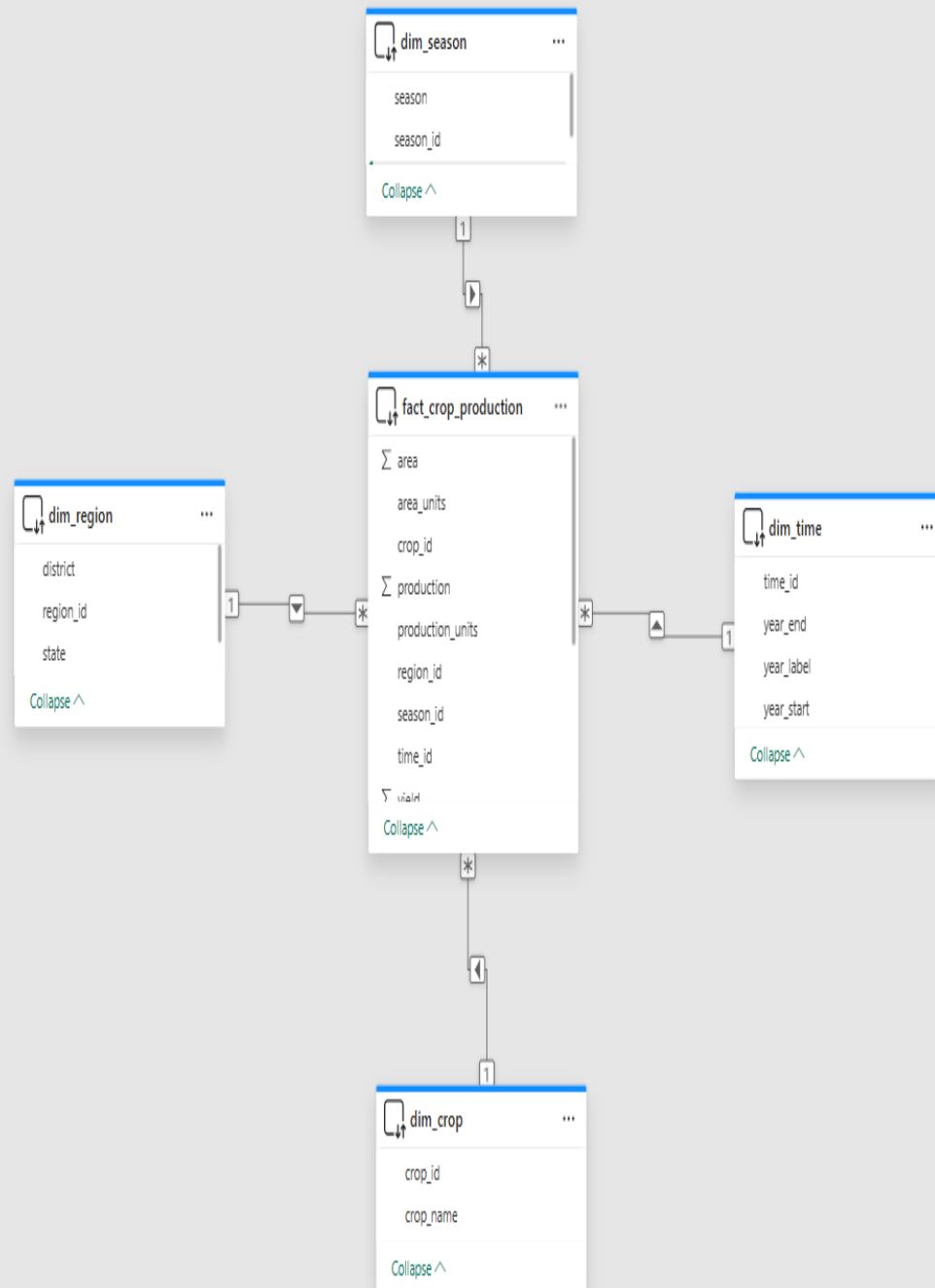# Data Modeling: Gold Layer Star Schema

## Dimension Tables

- **Crop:** Crop types and categories
- **Region:** State and district hierarchy
- **Season:** Agricultural season definitions
- **Time:** DateTime dimensions

## Why Star Schema?

- Simplifies complex analytics queries
- Optimizes Power BI performance
- Enables intuitive slicing, filtering, and drill-through navigation

📝 **Central Fact Table:** Crop Production metrics connecting all dimensions for comprehensive analysis

# Power BI: Executive Overview Dashboard

## High-Level KPIs

Total Production, Cultivated Area, and Average Yield metrics

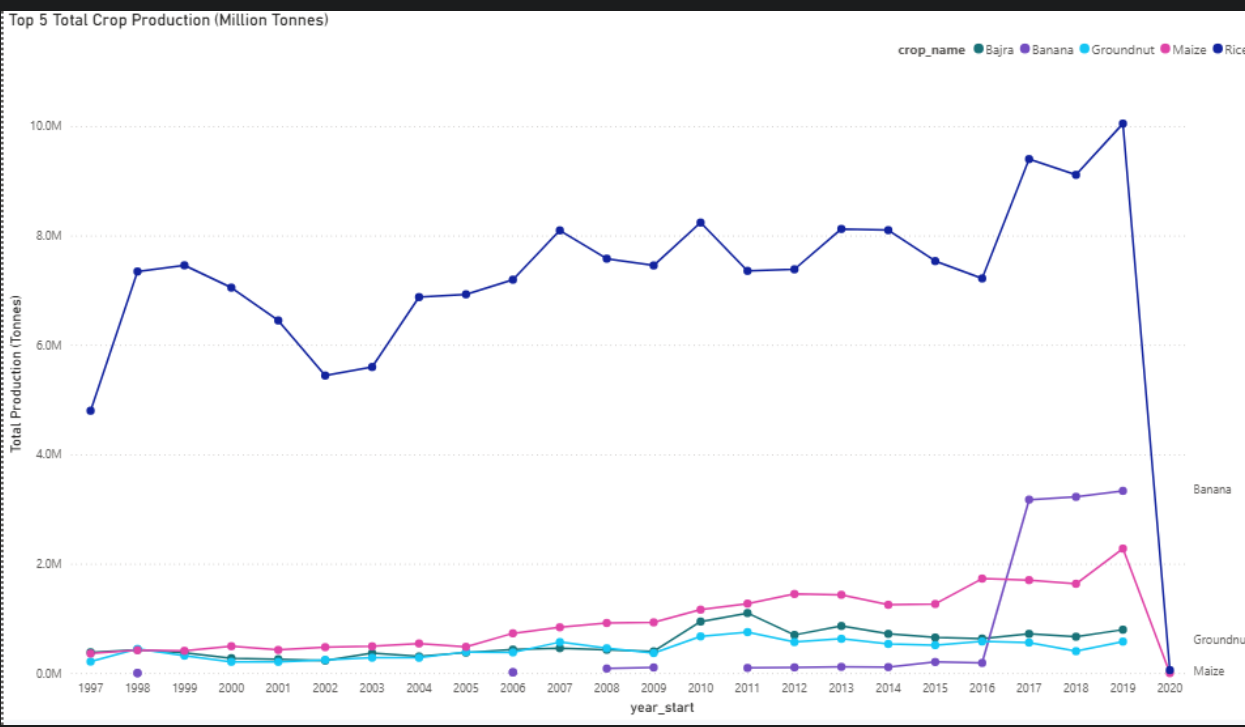## Interactive Controls
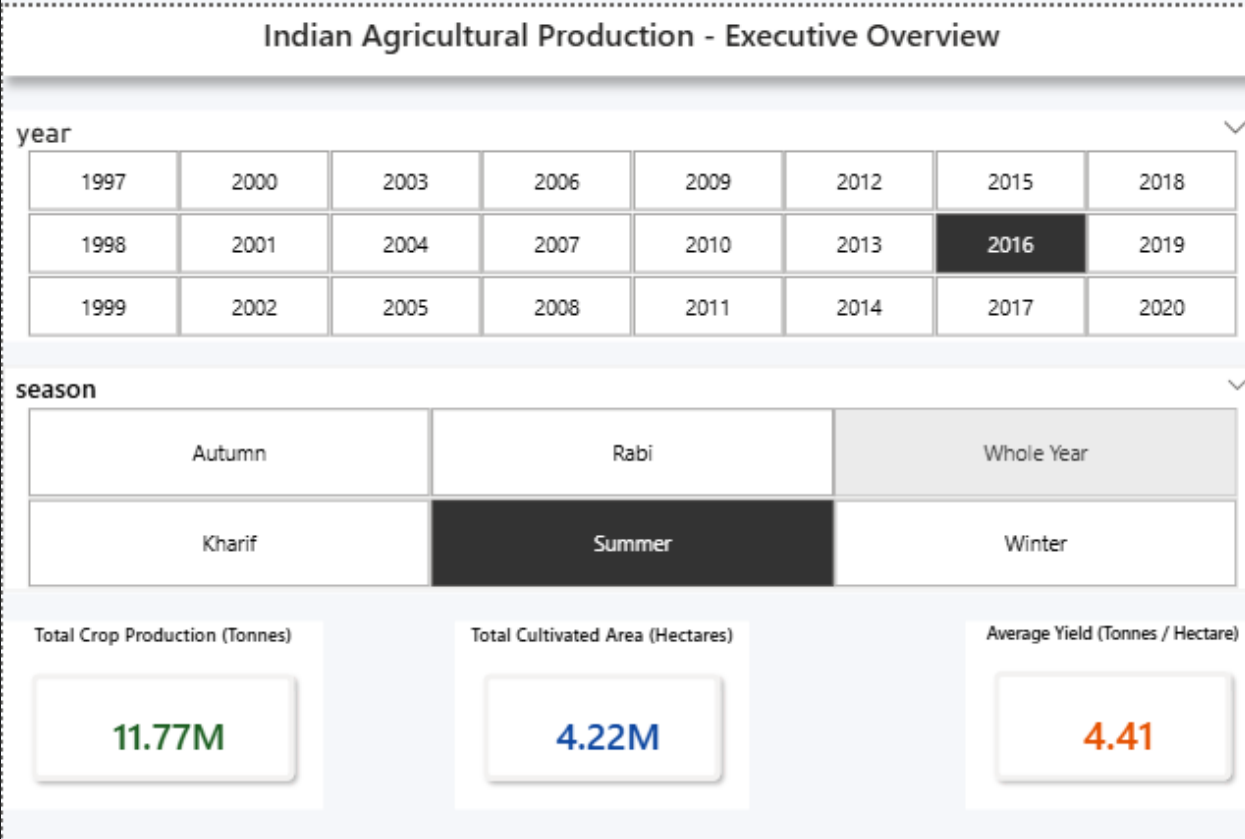
Year slicer for temporal analysis

## Top Performers
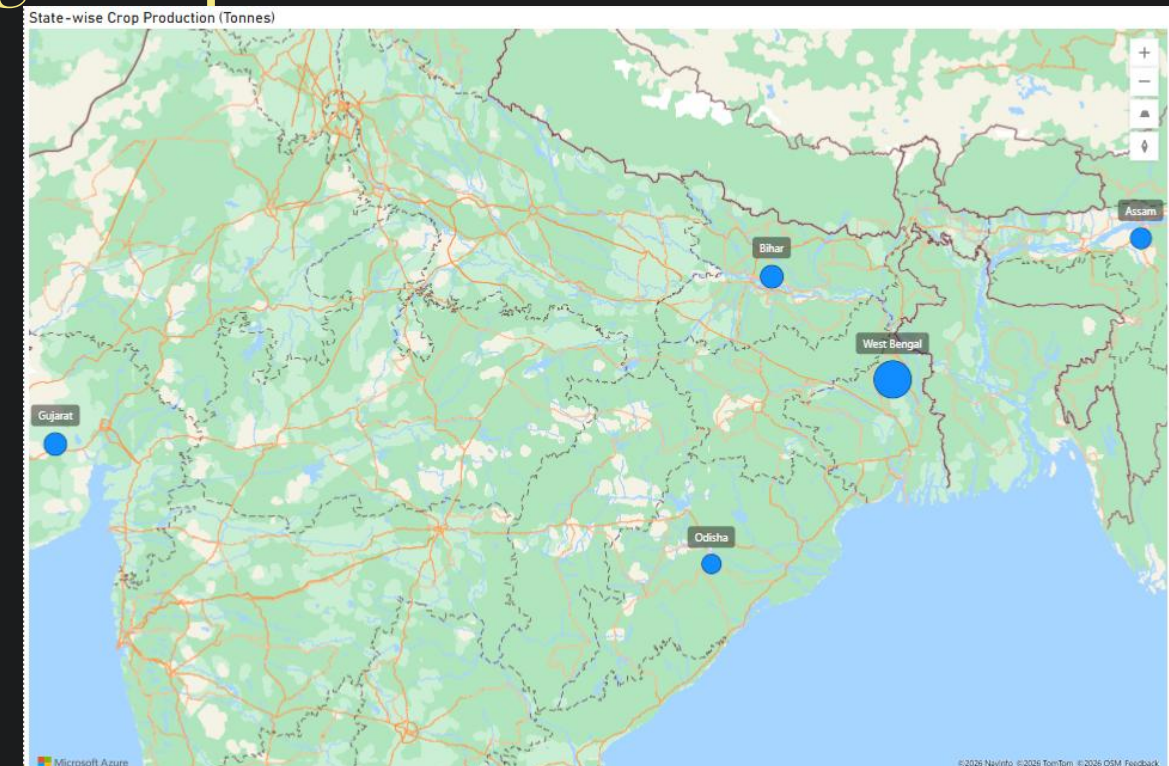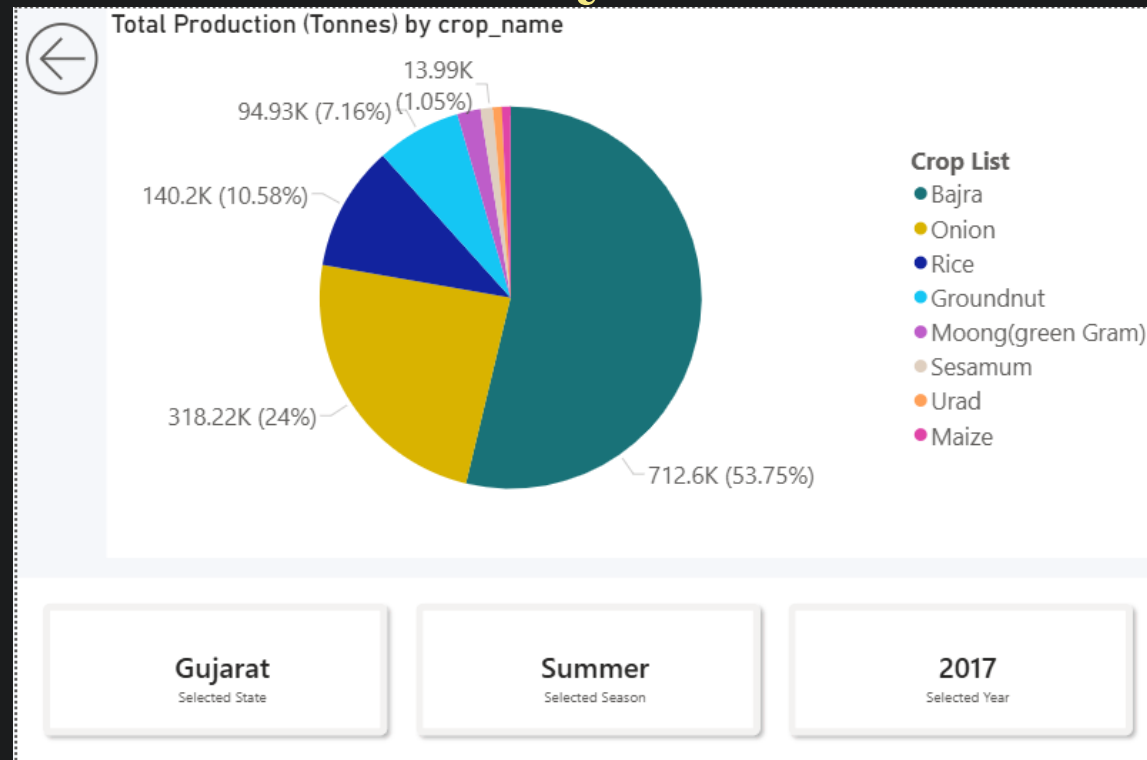
Leading crops by production volume

## Geographic Visualization

State-level production heat map

**Technical Integration:** Connected directly to Gold layer with optimized DAX measures for real-time performance

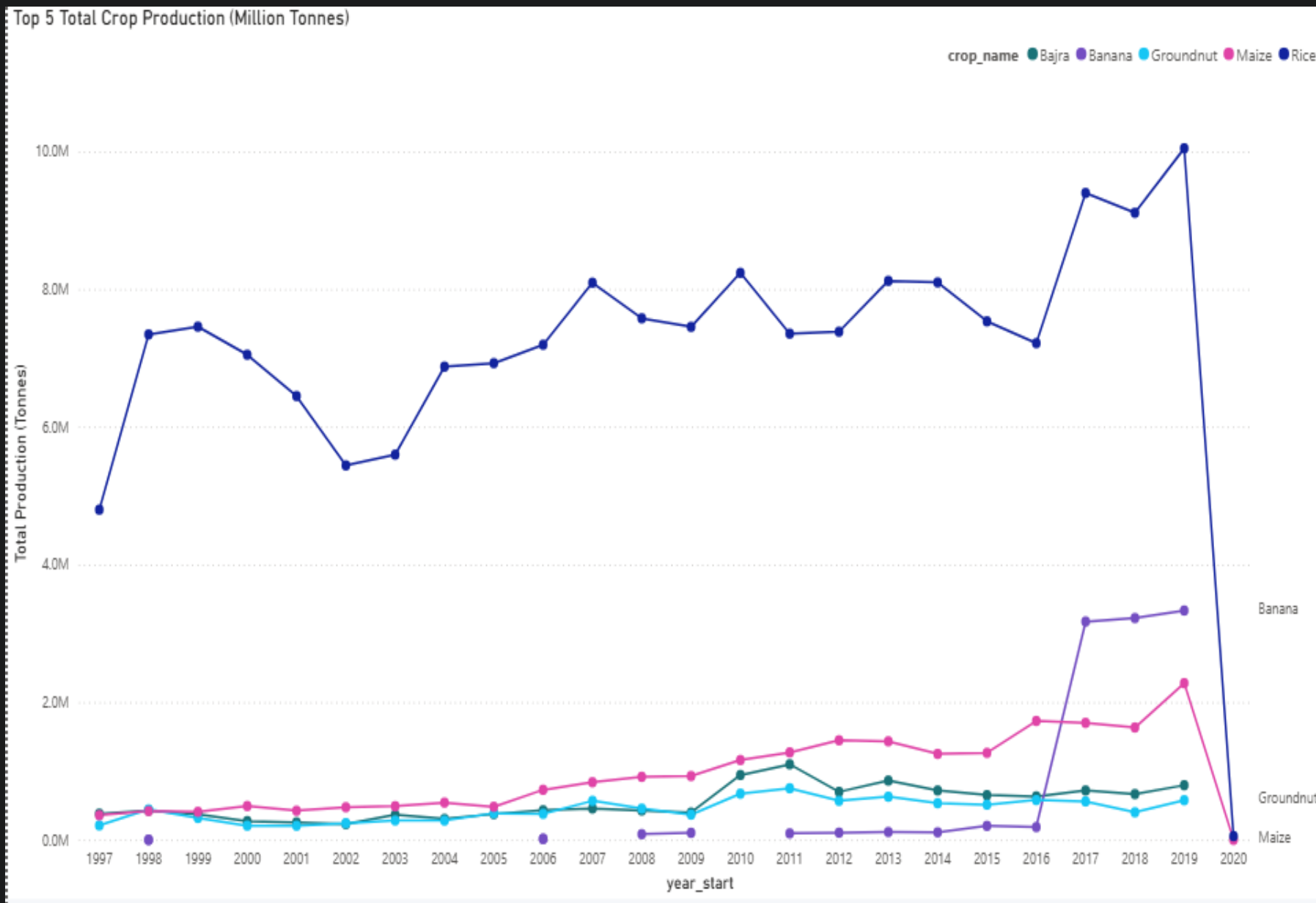# Advanced Analytics & Drill-Through Capabilities



## Enhanced Features

- **State-Level Drill-Through:** Navigate from overview to detailed state analysis

- **Dynamic KPI Cards:** Context-aware metrics for State, Season, and Year

- **Crop Production Breakdown:** Granular analysis by crop type

- **Multi-Year Trends:** Temporal patterns and growth trajectories

## Technical Highlights

- DAX with `SELECTEDVALUE` for context-sensitive calculations

- `REMOVEFILTERS` to control filter context propagation

- Year-independent trend visualizations

- Cross-filtering across multiple dimensions

# Analysis & Observation Achieved


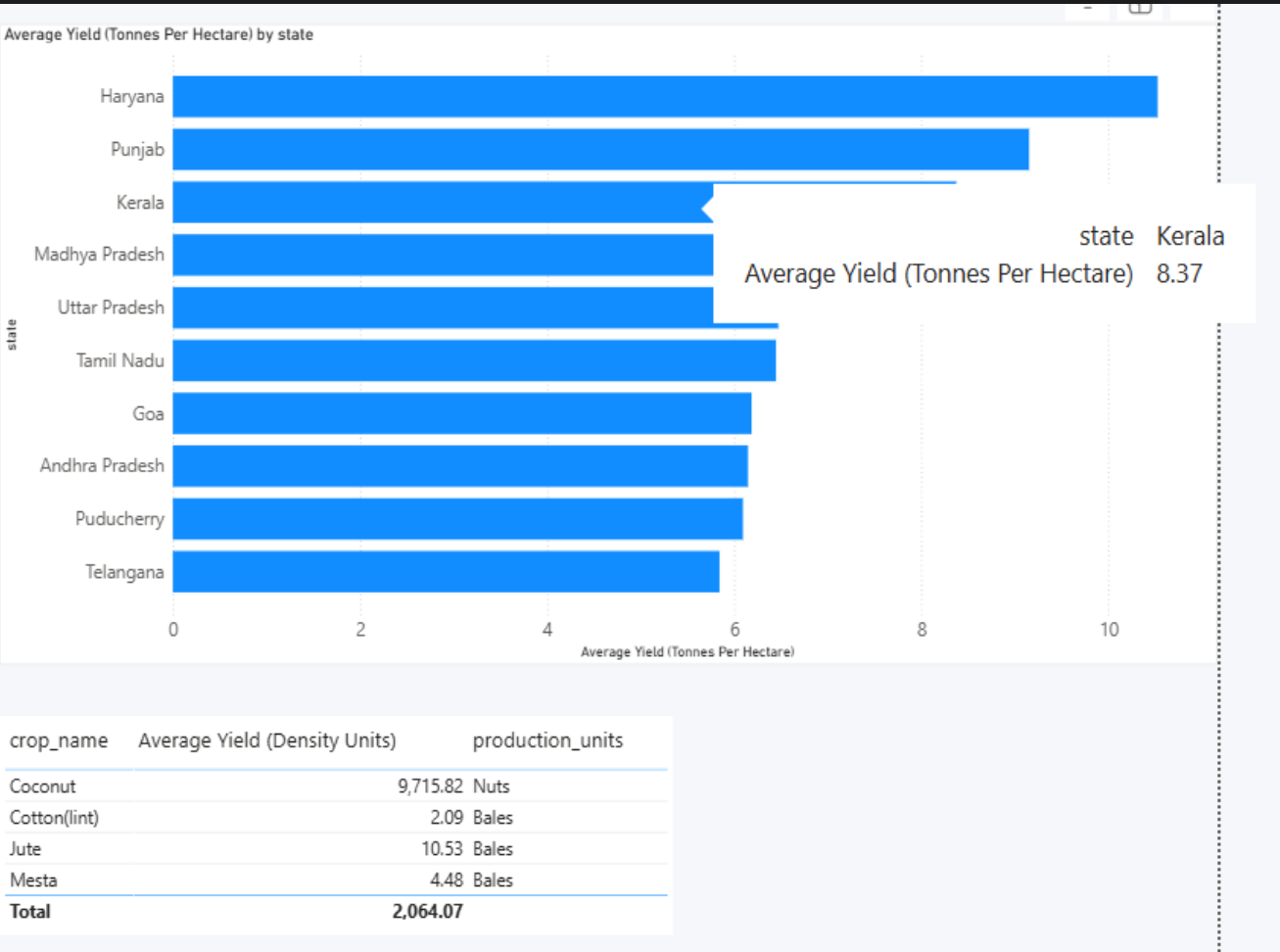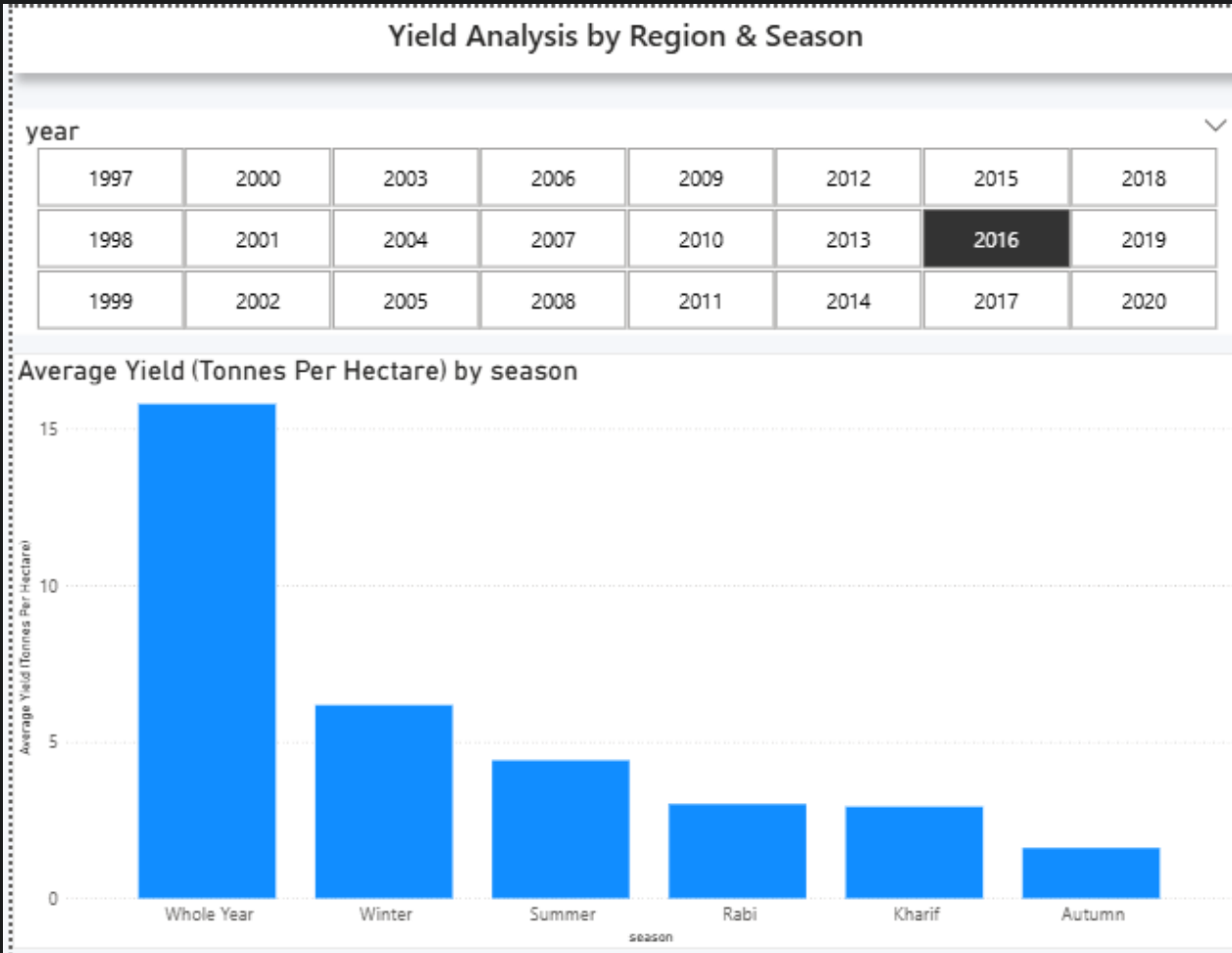Top 5 Total Crop Production (Million Tonnes)

**Business Insight**
Identifies:
- Crops with stable growth
- Crops with volatile production
- Long-term production shifts

Helps policymakers understand **which crops are becoming dominant or declining**

# Analysis & Observation Achieved



**Business Insight**

Reveals:

Which seasons are most productive

Which states have higher land productivity

# Key Technical Learnings & Outcomes

## Engineering Depth

- Spark vs Pandas trade-offs for different data volumes
- Handling inconsistent units across datasets
- Complex yield recalculation logic

## BI Optimization

- Managing filter context in Power BI
- DAX performance tuning
- Star schema modeling benefits

## Pipeline Design

- Modular ETL architecture with Airflow
- Separation of concerns across layers
- Scalable, maintainable design patterns

## Capstone Success

Delivered a fully automated, end-to-end analytics pipeline demonstrating enterprise-grade data engineering and BI implementation. All requirements successfully fulfilled with production-ready, scalable architecture.

"This project demonstrates real-world data engineering and analytics practices at scale."