

**CPSC 531-02**

**ADVANCED DATABASE MANAGEMENT**

**PROJECT  
CROP PRODUCTION METRICS**

**TEAM MEMBERS  
ANURAG GANJI (885174821)  
MIKAELLA SHARMA ( 885175406 )**

# **INTRODUCTION**

Crop production statistics give vital insights into the agricultural industry, allowing policymakers, academics, and farmers to analyze productivity, spot trends, and make educated decisions. Key criteria include yield per acre or hectare, production volume, crop diversification, and technological adoption. These measures aid in monitoring food security, assessing the impact of climate change, and directing resource allocation. They also help worldwide efforts toward sustainable agriculture and effective resource management.

Crop production metrics are important in formulating agricultural policies and plans in India. With a diversified agricultural terrain and a big population reliant on farming, indicators like overall production, crop-specific yields, and area are critical. Furthermore, examining elements like irrigation systems, fertilizer usage, and crop rotation patterns can aid in optimizing productivity, improving food security, and promoting sustainable practices. Crop production data that are accurate and timely enable stakeholders to solve difficulties and maximize India's agricultural potential, supporting rural development and economic progress.

# **PROJECT ABSTRACT**

The project's goal is to create a consolidated database system for crop production management in India. It will provide detailed information on crop varieties, cultivated land area, and overall yield. The database's goal is to provide significant information regarding crop production patterns in India, which will assist farmers, policymakers, and agricultural researchers. The project will leverage a variety of database systems to assure accuracy, efficiency, and usability. This program has enormous potential for aiding

informed decision-making, optimizing agricultural practices, and promoting the general growth and development of India's agriculture business.

## TECHNOLOGY USED

The project will make use of both the Windows and Mac operating systems. For efficient data processing and analysis, it will be built with Google Cloud Platform (GCP), Spark, and BigQuery. Python will be the major programming language, including important libraries such as PySpark, SparkML, and Pandas. The integrated development environment (IDE) will be Jupyter Notebook, which will provide a user-friendly and collaborative platform for coding and experimenting.

## FUNCTIONALITIES

### 1) Dataset Overview

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 345336 entries, 0 to 345335
Data columns (total 8 columns):
 #   Column      Non-Null Count  Dtype  
---  -- 
 0   State        345336 non-null   object 
 1   District     345336 non-null   object 
 2   Crop         345327 non-null   object 
 3   Crop_Year    345336 non-null   int64  
 4   Season       345336 non-null   object 
 5   Area          345336 non-null   float64
 6   Production   340388 non-null   float64
 7   Yield         345336 non-null   float64
dtypes: float64(3), int64(1), object(4)
memory usage: 21.1+ MB
```

```
print(dataset['crop'].unique())
['Areca nut' 'Arhar/Tur' 'Banana' 'Black pepper' 'Cashewnut' 'Coconut'
 'Cowpea(Lobia)' 'Dry chillies' 'Ginger' 'Groundnut' 'Maize'
 'Moong(Green Gram)' 'Oilseeds total' 'Other Kharif pulses'
 'other oilseeds' 'Rapeseed & Mustard' 'Rice' 'Sesamum' 'Sugarcane'
 'Sunflower' 'Sweet potato' 'Tapioca' 'Turmeric' 'Urad' 'Bajra'
 'Castor seed' 'Coriander' 'Cotton(lint)' 'Garlic' 'Gram' 'Guar seed'
 'Horse-gram' 'Jowar' 'Linseed' 'Masoor' 'Mesta' 'Niger seed' 'Onion'
 'Other Rabi pulses' 'Potato' 'Ragi' 'Safflower' 'Sannhamp'
 'Small millets' 'Soyabean' 'Tobacco' 'Wheat' 'Peas & beans (Pulses)'
 'Jute' 'Barley' 'Khesari' 'Moth' 'Other Cereals' 'Cardamom'
 'Other Summer Pulses' nan]
```

Figure: Dataset Overview

The project entails gathering data from Kaggle spanning the years 1997 to 2020 from various Indian states. The dataset is made up of eight instances (data columns include state\_name, district\_name, crop\_year, season, crop\_name, area, production, and yield). The data covers many agricultural seasons, including Autumn and Rabi. The primary goal is to forecast agricultural production, which is quantified in tonnes per hectare. The study hopes to obtain insights into agricultural production trends, identify variables impacting yield, and build prediction models to anticipate future crop yields by examining this dataset.

The code `dataset.info()` returns a summary of the dataset's information, such as the number of rows, columns, and data types for each column. It is used to acquire a brief overview of the structure and features of the dataset.

`dataset.info()` returns information such as :

1. The dataset's total number of entries (rows).
2. The number of columns and their individual names.
3. The number of non-zero values in each column.
4. Each column's data type, such as integer, float, text, or datetime.
5. Additional memory utilization data.

The dataset of code `snippets['Crop'].unique()` is used to retrieve the unique values from the dataset's 'Crop' column. For example, because the dataset contains information about various crops such as wheat, rice, corn, and barley, it is referred to as the `dataset['Crop']`. The `unique()` method will return an array or list of unique crop names, such as ['Wheat', 'Rice', 'Corn', 'Barley']. This function assists in identifying all of the unique crop kinds included in the dataset, allowing for additional analysis or filtering based on crop classifications.

## **2) Data Cleaning**

In the context of the "Crop Production Statistic" dataset, data cleaning refers to the process of finding and resolving flaws in the dataset in order to enhance its quality, consistency, and dependability. Data cleaning procedures can be used to manage missing values, duplicates, guarantee uniform formatting, discover and handle outliers, fix data types, resolve inconsistent values, and validate data.

The following are possible data cleaning processes that are done on the project dataset:

### **1. Data Imputation - Missing Values**

The initial stage in pre-processing is to search the dataset for missing values. Handling missing values in the project dataset entails identifying columns where significant information, such as crop yield or area, is missing. One method was to eliminate the rows with missing values, which ensures that the remaining data is complete while potentially lowering the total dataset size. Missing values are imputed using methods such as mean or median imputation, which utilize the average or middle value of the available data to fill in the gaps. Imputation allows valuable data to be retained while assuring a complete dataset for study. The decision between removal and imputation is influenced by the extent and significance of missing values, as well as the possible impact on future analysis and interpretations.

### **2. Removing Duplicates**

We thoroughly checked the dataset for duplicate items. To assure data integrity, we examined crop, state, and year combinations in particular. During our examination, we discovered and eliminated any duplicate records. We guaranteed that each unique combination of crop, state, and year was represented only once in the dataset by

removing these duplicates. Duplicate entries might add biases and skew the analytic results, thus this step was critical to maintaining the quality and dependability of our data. We made certain that the dataset was clear of duplicate information, allowing for more accurate and useful insights regarding agricultural output in various Indian states over time.

### **3. Consistent Formatting**

We understood the significance of uniform formatting in the dataset. We concentrated on standardizing the formatting of columns such as crop names, state names, and years to accomplish this. We went through each entry carefully, making sure that any discrepancies in capitalization or abbreviations were rectified. We eradicated discrepancies and established uniformity across the dataset by standardizing the formatting. This step was critical in improving data clarity and facilitating proper analysis. We meticulously made the necessary formatting adjustments while working together, resulting in a dataset with consistent and well-formatted entries. This resulted in smoother data interpretation and the ability to derive useful insights from the information.

### **4. Outlier Detection and Handling**

We recognized the need of finding and dealing with outliers in the dataset. We worked to find possible outliers and assess their validity using indicators such as yield per acre. We determined whether these extreme numbers were legitimate data points or erroneous entries using rigorous analysis. Where applicable, we handled outliers by deleting them or applying appropriate adjustments to lessen their influence on analytic results. Our collaborative work guaranteed that the dataset was free of outlier distortions, enabling for more accurate and dependable insights to be generated from the data.

## **5. Data Validation and Type Correction**

We prioritized data validation and assuring the dataset's integrity. To ensure accuracy and consistency, we completed data validation checks together. For example, we ensured that crop yields were within realistic limits and cross-validated data across relevant columns to find any anomalies. Furthermore, we thoroughly evaluated the data types of columns to ensure that numerical values were recorded as numeric data types rather than text or category variables. Our collaborative effort in executing these checks and altering data types as needed enhanced the dataset's dependability, improving the accuracy and quality of future analyses and interpretations.

When these data cleaning procedures are applied to the "Crop Production Statistics" dataset, the resultant cleaned dataset will have higher quality, consistency, and dependability. This allows for more precise agricultural production metrics research and insights in India.

### **3) Data Analysis**

In our research, we used modern technologies like SparkML and PySpark to create reliable crop production prediction models. We processed and analyzed massive amounts of data from the dataset effectively by harnessing the power of distributed computing provided by Spark. Using historical data on crop yield, area, and other pertinent factors, we constructed strong prediction models. We will construct multiple regression techniques using PySpark machine learning modules. These models accurately predicted crop yield by capturing the intricate interactions between input factors and crop yield.

We processed the dataset effectively by using Spark's scalability and the distributed computing environment of Google Cloud Platform (GCP), resulting in speedier model

training and assessment. BigQuery was used for data preparation and feature engineering activities, taking advantage of its high-performance storage and querying capabilities.

We used Jupyter Notebook as the IDE for our investigation, which provides an interactive and collaborative environment for code creation, model experimentation, and visualization. Panda libraries were helpful in studying the dataset, displaying patterns, and learning about agricultural production trends. We constructed sophisticated prediction models that used previous data to reliably project future crop yields. This improves decision-making for farmers, policymakers, and agricultural stakeholders, allowing for better planning, resource allocation, and production. This combination of technology and libraries enabled us to get useful insights, reveal trends, and promote future machine learning processes for improved crop production decision-making.

## 1. Groupby Analysis

```
In [8]:  
    crop['Season'].unique()  
  
Out[8]:  
    array(['Kharif      ', 'Whole Year ', 'Autumn      ', 'Rabi      ',  
          'Summer      ', 'Winter      '], dtype=object)
```

Figure: Groupby Season data

### a. Groupby Seasons

We performed group-by analysis on the dataset by season and crop to acquire insight into agricultural production trends and variances. We were able to investigate the trends and features peculiar to each season by categorizing the data according to seasons such as Autumn, Rabi, and others. We pooled the data for each season to generate key parameters such as average crop yield, total area under cultivation, and production amount. This

research enabled us to comprehend the seasonal fluctuations in agricultural yield. We found the seasons with the highest and lowest crop yields, providing farmers and policymakers with useful information for optimizing farming techniques and resource allocation depending on seasonal trends.

### **b. Groupby Crops**

Similarly, we conducted crop-specific group-by analysis, aggregating the data to analyze crop-specific trends and performance. This investigation revealed information about the most regularly grown crops, their average yields, and production fluctuations throughout time. We discovered crops with high yields and consistent production patterns, allowing stakeholders to make educated decisions about crop selection, market demand, and investment prospects. We were able to get useful insights into the dynamics of crop production in respect to seasonal fluctuations and particular crop selections by exploiting the dataset's information and doing group-by analysis by season and crop. These insights can help with better decision-making, better agricultural practices, and overall agricultural growth.

## **4 ) Google Cloud Platform**

In our project, we used Google Cloud Platform (GCP) and, in particular, Cloud Dataproc to improve our data analysis on the dataset. We got access to a strong and scalable infrastructure through GCP, which permitted the efficient processing and analysis of agricultural production data. We built clusters to manage massive amounts of data using Cloud Dataproc, including a single-node and a three-node cluster. These clusters have Apache Spark installed, allowing us to undertake distributed data processing and analysis. We obtained faster execution and better performance by dividing the workload over numerous nodes.

To store and maintain our dataset, we used Cloud Storage, a scalable and long-lasting object storage solution. This allowed for smooth interaction with Cloud Dataproc and guaranteed quick access to data during analysis. We were able to properly process and analyze the dataset using Cloud Dataproc and Apache Spark's distributed computing capabilities. We managed to do sophisticated computations and gain insights from agricultural production data because of Spark's parallel processing capabilities.

We improved our data analysis workflow by integrating GCP and Cloud Dataproc, allowing for quicker execution, scalability, and effective computing resource use. This enabled us to identify useful insights and trends in crop production data, allowing farmers, policymakers, and stakeholders in the agriculture sector to make better-informed decisions.

#### **4 ) Data Exploration and Visualization**

We used the Google Cloud Platform (GCP) in our project to help with data visualization for agricultural production statistics.

This is how we executed it:

1. We uploaded the crop production dataset to Google Cloud Storage, which functioned as our primary data storage option. This made the dataset more accessible and manageable.
  
2. We used GCP's Dataproc service to construct our Spark project. Within Dataproc, we set up a cluster with one master node. This networked computing architecture allowed for the efficient processing of massive amounts of data.

3. We constructed a PySpark Jupyter notebook within the cluster to provide an interactive and collaborative environment for data processing and visualization. We used Dataproc to read the dataset from cloud storage, which easily connected cloud storage with the Hadoop Distributed File System (HDFS).
4. We ran queries on the dataset using PySpark SQL to extract important crop production information. We then stored the query results to the cloud for convenient access and subsequent analysis.

We were able to study agricultural production figures and get significant insights by utilizing these strategies. The usage of GCP, Dataproc, and PySpark SQL created a scalable and efficient data visualization platform, allowing us to efficiently display and discuss crop production trends, patterns, and statistics. This visualization component improved our comprehension of the dataset and aided farmers in the agriculture sector in making informed choices.

In addition after doing data analysis on the agricultural production dataset, we saved the results in Google BigQuery tables. This stage aided in the creation of a direct link between the examined data and Tableau Online, a sophisticated and comprehensive data visualization tool.

We were able to swiftly and effectively visualize crop production facts and trends by exploiting this relationship. Tableau Online provided a plethora of dynamic and visually attractive choices for producing informative charts, graphs, and interactive dashboards based on BigQuery data.

This connection of Google BigQuery and Tableau Online bridges the gap between data analysis and visualization, allowing us to display the agricultural production dataset's major results in an interesting and instructive manner. The combination of robust data

analytic tools and clear and aesthetically attractive visualizations improved our ability to share and grasp the dataset's important conclusions.

#### **4 ) BigQuery**

BigQuery proven to be an excellent tool for effectively handling and querying enormous amounts of data. We used BigQuery to store and organize the dataset, allowing for easy access and retrieval of pertinent data. We were able to extract particular subsets of data, apply aggregations, and run extensive analysis on crop production parameters using its high-performance querying capabilities. Furthermore, the integration of BigQuery and other technologies, such as Tableau Online, enabled us to link and view the dataset directly, resulting in dynamic and meaningful visual representations of crop production trends and patterns. Google BigQuery's scalability and versatility presented us with a strong tool.

We were able to find interesting patterns, trends, and correlations within the crop production information by harnessing the capabilities of Google BigQuery, allowing us to make data-driven choices and apply successful methods to improve agricultural productivity and sustainability.

# ARCHITECTURE AND DESIGN

We utilized Spark and PySpark to clean and analyze the dataset in our agricultural production study project, assuring reliable and consistent data for subsequent research.

To begin, we installed Spark and created a Spark Session, which serves as an entry point for Spark programming with the Dataset and Data Frame API. This session taught us how to build Data Frames, register them as tables, and perform SQL queries on the data. We successfully generated a Spark session by using the `getOrCreate()` function and the builder pattern method `builder()`. We next loaded our CSV-formatted dataset into a Data Frame for additional processing.

Following that, we used PySpark to cleanse the data, which was an important step in ensuring the dataset's correctness and consistency. We installed PySpark using the PIP Install PySpark command and utilized its capabilities to eliminate duplicates, errors, and unnecessary data from our dataset. We were able to effectively cleanse data by creating Data Frames, registering them as tables, and performing SQL queries. With a clean dataset, we examined the data using several ways, including building data frames and assessing the data using SQL queries.

We placed the cleansed dataset into a Google Cloud Storage bucket for analysis. This enabled us to securely store and retrieve the dataset, allowing for efficient data analysis through cloud-based technology. We cleaned and analyzed the agricultural production dataset using the capabilities of Spark and PySpark. These technologies gave us the tools and features we needed to verify data quality, minimize discrepancies, and get significant insights into crop production patterns. The combination of Spark, PySpark, and cloud storage allows us to do comprehensive data analysis in a scalable and fast way, boosting agricultural decision-making processes.

The screenshot shows the 'Create a bucket' wizard in the Google Cloud Platform. The left sidebar includes 'Buckets', 'Monitoring NEW', and 'Settings'. The main area has four steps completed: 'Name your bucket' (Name: crop\_dbms), 'Choose where to store your data' (Location: us (multiple regions in United States), Location type: Multi-region), 'Choose a storage class for your data' (Default storage class: Standard), and 'Choose how to control access to objects' (Public access prevention: On, Access control: Uniform). Step five, 'Choose how to protect object data', is currently selected. It contains three protection tools: 'None' (selected), 'Object versioning (for data recovery)' (disabled), and 'Retention policy (for compliance)' (disabled). A 'Good to know' sidebar provides location pricing information and current configuration details. A cost estimation table shows rates for standard storage.

Figure: Bucket Creation

The screenshot shows the 'Bucket details' page for 'crop\_dbms'. The left sidebar includes 'Marketplace' and 'Release Notes'. The main area displays bucket metadata: Location (us (multiple regions in United States)), Storage class (Standard), Public access (Not public), and Protection (None). Below this are tabs for 'OBJECTS', 'CONFIGURATION', 'PERMISSIONS', 'PROTECTION', 'LIFECYCLE', 'OBSERVABILITY NEW', and 'INVENTORY REPORTS NEW'. The 'OBJECTS' tab shows a single file, 'AgricultureDataset.csv', uploaded successfully. A message at the bottom indicates '1 file successfully uploaded'. A status bar at the bottom right shows 'Uploads and First Project dbms operations' with a green checkmark for 'AgricultureDataset.csv'.

Create table

### Source

Create table from — Google Cloud Storage

Select file from GCS bucket or use a URI pattern \* —  crop\_dbms/AgricultureDataset.csv [BROWSE](#) [?](#)

File format — CSV

Source Data Partitioning

### Destination

Project \* — first-project-dbms [BROWSE](#)

Dataset \* — crop\_dataset

Table \* — Crop\_table  
Unicode letters, marks, numbers, connectors, dashes or spaces allowed.

Table type — Native table

### Schema

Auto detect

[CREATE TABLE](#) [CANCEL](#)

Figure: Creating relational table

Free trial status: \$300.00 credit and 91 days remaining - with a full account, you'll get unlimited access to all of Google Cloud Platform. [DISMISS](#) [ACTIVATE](#)

Google Cloud First Project dbms clusters [X](#) [Search](#) [M](#)

**Dataproc** [Create a Dataproc cluster on Compute Engine](#)

**Jobs on Clusters** [Clusters](#) [Jobs](#) [Workflows](#) [Autoscaling policies](#)

**Serverless** [Batches](#)

**Metastore Services** [Metastore](#) [Federation](#)

**Utilities** [Component exchange](#) [Workbench](#) [Release Notes](#)

**Set up cluster** Begin by providing basic information.

**Configure nodes (optional)** Change node compute and storage capabilities.

**Customize cluster (optional)** Add cluster properties, features, and actions.

**Manage security (optional)** Change access, encryption, and security settings.

[CREATE](#) [CANCEL](#)

[EQUIVALENT COMMAND LINE](#)

**Components**

**Component Gateway**

**Enable component gateway** Provides access to the web interfaces of default and selected optional components on the cluster. [Learn more](#)

**Optional components**

Anaconda [?](#)  
 Hive WebHCat [?](#)  
 Jupyter Notebook [?](#)  
 Zeppelin Notebook [?](#)  
 Druid [?](#)  
 Presto [?](#)  
 ZooKeeper [?](#)  
 Ranger [?](#)  
 HBase [?](#)  
 Flink [?](#)  
 Docker [?](#)  
 Solr [?](#)  
 Hudi [?](#)

Free trial status: \$300.00 credit and 91 days remaining - with a full account, you'll get unlimited access to all of Google Cloud Platform.

DISMISS ACTIVATE

Google Cloud First Project dbms clusters X Search M

## Dataproc

### Create a Dataproc cluster on Compute Engine

**Jobs on Clusters**

- Clusters
- Jobs
- Workflows
- Autoscaling policies

**Serverless**

- Batches

**Metastore Services**

- Metastore
- Federation

**Utilities**

- Component exchange
- Workbench
- Release Notes

**Initialization actions**

Set up cluster  
Configure nodes (optional)  
Customize cluster (optional)  
Manage security (optional)

+ ADD PROPERTIES

bucket/folder/file 1 —  dataproc-staging-us-central1-1026196437002-eilfhbjb googl BROWSE

+ ADD INITIALIZATION ACTION

**Custom cluster metadata**

Add custom metadata to cluster instances. [Learn more](#)

Key 1 *	spark-bigquery-connector-version	Value 1	0.21.0
---------	----------------------------------	---------	--------

+ ADD METADATA

**Scheduled deletion**

Use Scheduled Deletion to help avoid incurring Google Cloud charges for an inactive cluster. [Learn more](#)

Delete on a fixed time schedule  
 Delete after a cluster idle time period without submitted jobs

CREATE CANCEL EQUIVALENT COMMAND LINE

Figure: Cluster creation and BigQuery key from DataProc

Free trial status: \$300.00 credit and 91 days remaining - with a full account, you'll get unlimited access to all of Google Cloud Platform.

DISMISS ACTIVATE

Google Cloud First Project dbms clusters X Search M

## Dataproc

### Clusters

**Jobs on Clusters**

- Clusters
- Jobs
- Workflows
- Autoscaling policies

**Serverless**

- Batches

**Metastore Services**

- Metastore
- Federation

**Utilities**

- Component exchange
- Workbench
- Release Notes

**PERMISSIONS** **LABELS**

No clusters selected

Please select at least one resource.

Request to create cluster cluster-6ddecrop submitted X

# DEPLOYMENT INSTRUCTIONS

In our project, we went through a number of processes to create the environment and tools needed for effective data analysis and visualization.

We began by downloading and installing Apache Spark from its official website, which served as the foundation for our data processing and analysis operations. This distributed computing platform enabled us to successfully manage massive amounts of data.

Following that, we set up Jupyter Notebook and installed the Python components needed to build an interactive and user-friendly writing environment. This interface allowed us to operate easily with Spark and Python, easing data handling and analysis.

We built up a cluster in the Google Cloud Platform (GCP) using Dataproc to take use of the capabilities of cloud computing. This gave us a scalable and efficient computer environment in which to process our agricultural production dataset.

We installed the essential libraries and dependencies to support our code in preparation for our data analysis. In addition, we constructed a Google Cloud Storage bucket to safely store our dataset.

Finally, we explored and gained insights from our agricultural production data using BigQuery, a sophisticated data analysis and visualization tool. We were able to conduct SQL queries, analyze data, and produce relevant visualizations by integrating BigQuery into our Jupyter Notebook environment.

We developed a robust framework for data analysis and visualization through these phases, allowing us to gain valuable insights and make acquired crop production decisions.

# STEPS TO RUN

**Step 1:** To begin our project, we logged into the Google Cloud Platform (GCP) application. We developed a new project especially customized to our crop production study after gaining access to our GCP account.

The screenshot shows the Google Cloud Compute Engine interface. The left sidebar includes sections for Compute Engine, Virtual machines (selected), Instance templates, Sole-tenant nodes, Machine images, TPUs, Committed use discounts, Reservations, Migrate to Virtual Machines, Storage (Disks, Snapshots, Marketplace), and Release Notes. The main content area displays 'VM instances' with tabs for Instances, Observability, and Instance Schedules. Under Instances, there is a table listing three VM instances: 'cluster-6ddecrop-m' (Status: Running, Name: cluster-6ddecrop-m, Zone: us-central1-f, Connect: SSH), 'cluster-6ddecrop-w-0' (Status: Running, Name: cluster-6ddecrop-w-0, Zone: us-central1-f, Connect: SSH), and 'cluster-6ddecrop-w-1' (Status: Running, Name: cluster-6ddecrop-w-1, Zone: us-central1-f, Connect: SSH). Below the table is a 'Related actions' section with six cards: Explore Backup and DR (NEW), View billing report, Monitor VMs, Explore VM logs, Set up firewall rules, Patch management, and Load balance between VMs.

Figure: Creating Virtual Machine instance on the Cloud Engine

After successfully creating the project, we proceeded to establish a virtual machine (VM) instance within the GCP architecture. The virtual machine instance acted as our computing environment, giving us the tools and capabilities we needed to complete our data processing activities. Various characteristics, such as machine type, disk size, and operating system, were defined during the VM instance formation procedure. These parameters were chosen based on our project's requirements and projected workload.

We received remote access to the VM instance after it was properly provisioned, allowing us to login and interact with it. We were able to administer the VM remotely, install software packages, and run commands for data processing and analysis.

**Step 2:** To guarantee adequate network connectivity and security, we modified the firewall settings and launched the VM instance from the Google Cloud Platform (GCP).

The screenshot shows the Google Cloud Platform interface for creating a new VM instance. The top navigation bar includes a free trial status message, a dismiss button, and an activate button. The main header shows "Google Cloud" and the project name "First Project dbms". A search bar and various navigation icons are also present.

The left sidebar lists options for creating a VM instance:

- New VM instance**: Create a single VM instance from scratch.
- New VM instance from template**: Create a single VM instance from an existing template.
- New VM instance from machine image**: Create a single VM instance from an existing machine image.
- Marketplace**: Deploy a ready-to-go solution onto a VM instance.

The main configuration area is titled "Create an instance". It includes a "Deploy a container image to this VM instance" section with a "DEPLOY CONTAINER" button, and a "Boot disk" section where the user has selected "crop-production1" (Type: New balanced persistent disk, Size: 10 GB, License type: Free, Image: Debian GNU/Linux 11 (bullseye)). There is also a "CHANGE" button for the boot disk settings.

The "Identity and API access" section includes "Service accounts" (set to "Compute Engine default service account") and "Access scopes" (radio buttons for "Allow default access" (selected), "Allow full access to all Cloud APIs", and "Set access for each API").

The "Pricing summary" section displays a monthly estimate of \$25.46, which is approximately \$0.03 hourly. It details the cost breakdown for 2 vCPU + 4 GB memory (\$24.46), 10 GB balanced persistent disk (\$1.00), and a total of \$25.46. A link to "Compute Engine pricing" and a "LESS" button are also shown.

Figure : Firewall Settings on VM

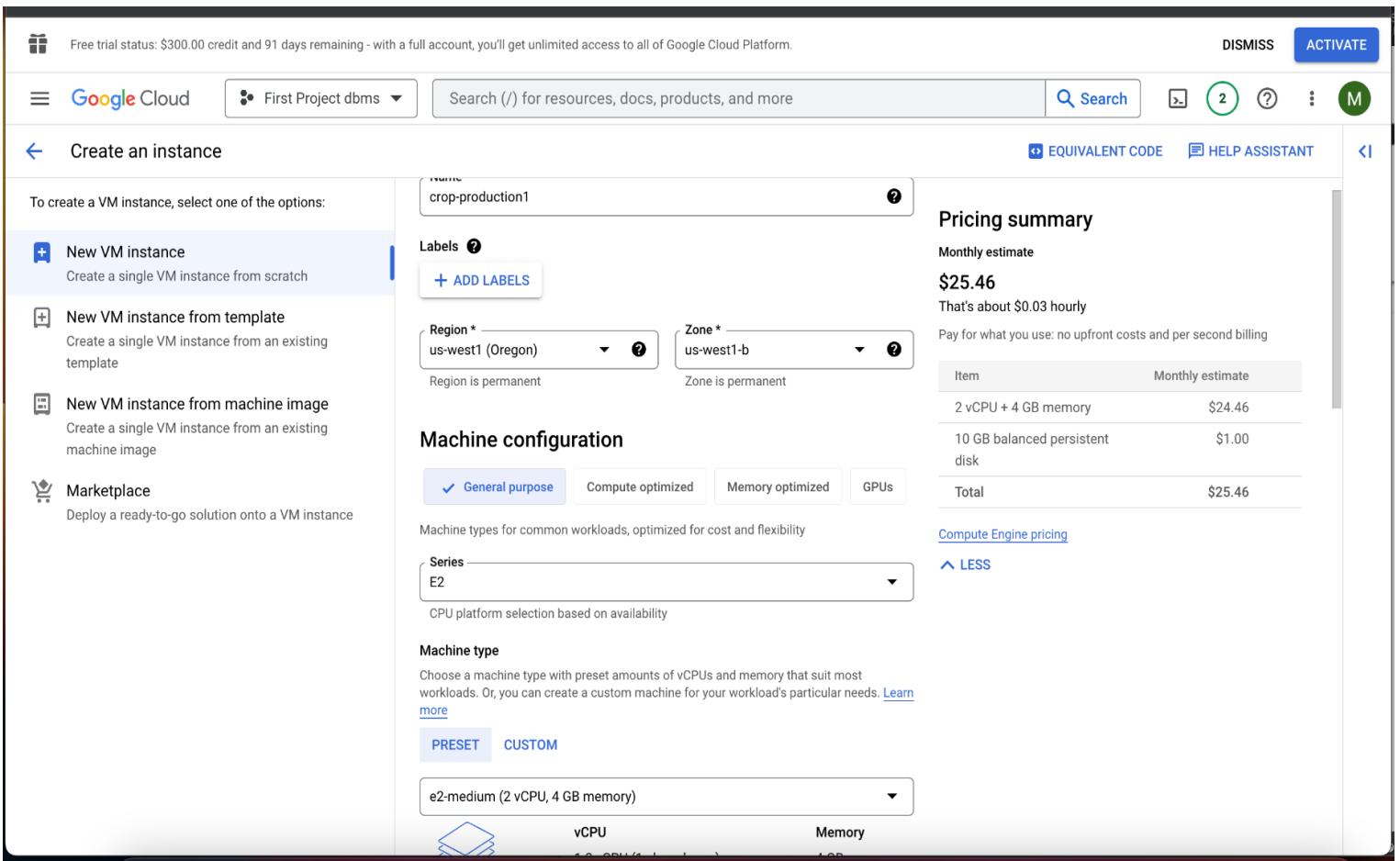


Figure : Machine Configuration on VM Instance

First, we went to the GCP interface and went to the networking area. We found the firewall rules in the networking setup and performed the appropriate changes. We launched the VM instance after upgrading the firewall rules. This started the boot-up process, making the virtual computer operational and ready to use. We kept an eye on the instance's status and confirmed its successful startup.

We assured that essential network access was permitted while preserving proper security measures by adjusting the firewall settings and launching the VM instance. This stage supplied us with a working computing environment in which to carry out our crop production analytic duties.

**Step 3:** Using an SSH terminal, we installed Jupyter Notebook and other essential programs to create the development environment for our project.

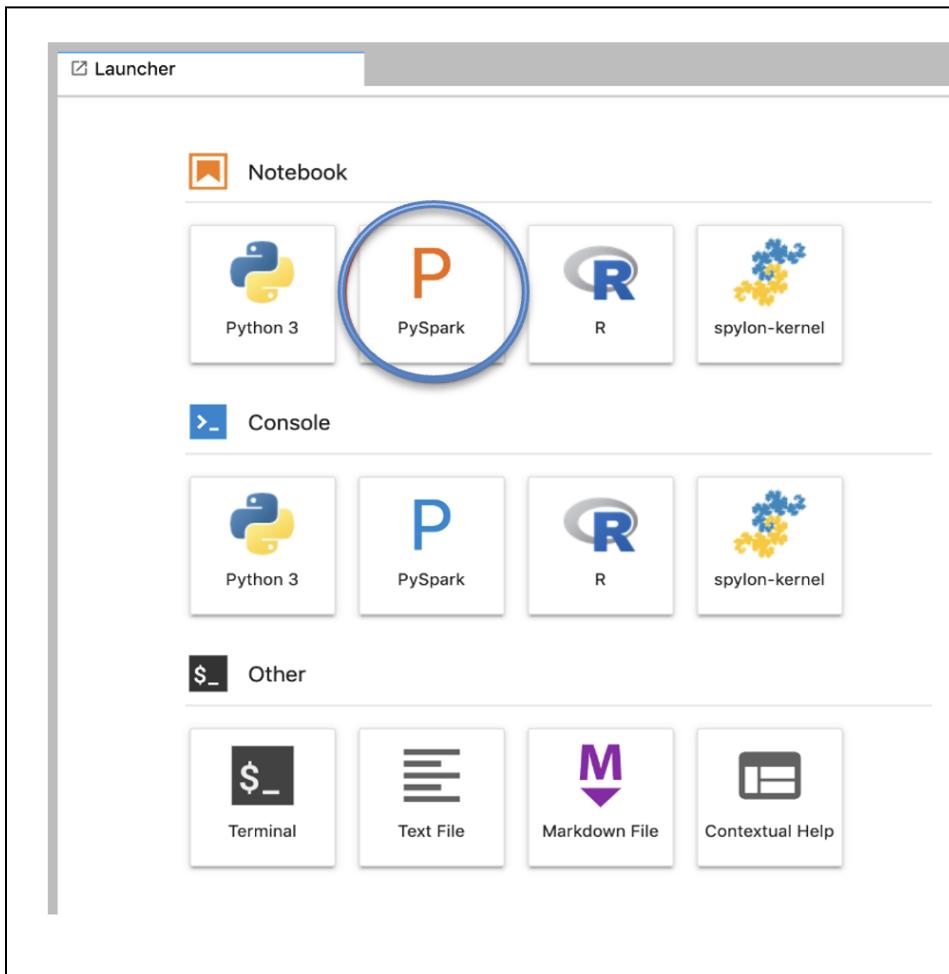


Fig: Launching Pyspark

After successfully connecting to the VM instance via SSH, we ran the necessary commands to install Jupyter Notebook. This included utilizing package managers like pip or conda to guarantee that the most recent version was installed. In addition, we installed additional necessary packages and libraries for our data analysis and processing operations. NumPy, Pandas, Matplotlib, Seaborn, and any additional dependencies relevant to our project requirements may be included in these packages.

We carefully followed the directions supplied by the package managers throughout the installation procedure to guarantee that the installs were successful. We created a stable development environment within the VM instance by installing Jupyter Notebook and other required programs. This environment enabled us to rapidly develop, execute, and collaborate on code using the Jupyter Notebook interface, allowing for seamless integration.

**Step 4:** After configuring the VM server, we launched Jupyter Notebook and configured it to run our data visualization scripts.

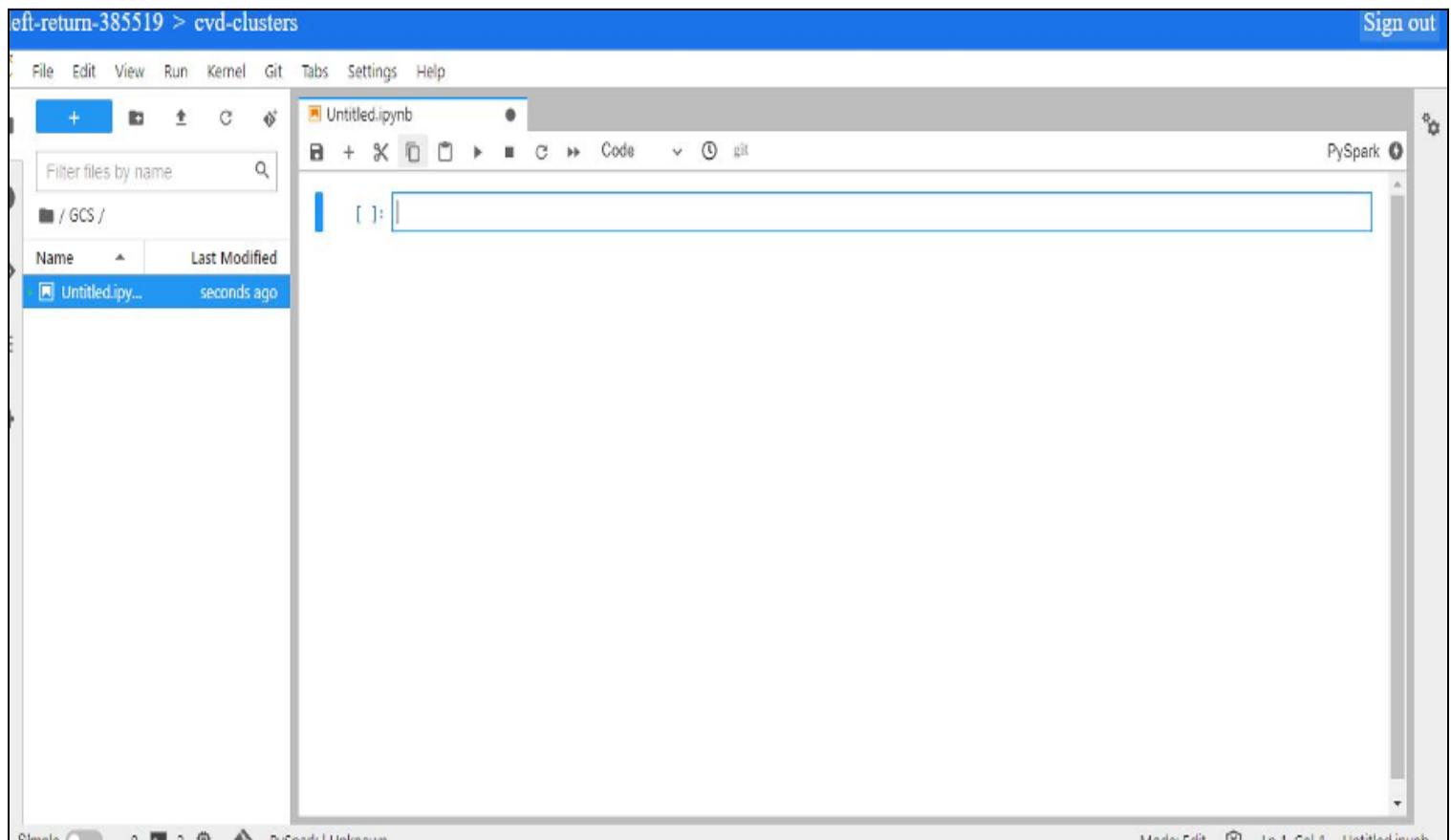


Fig: Launching Jupyter Notebook

To start the Jupyter Notebook server, we used the SSH terminal and entered the proper command, giving the desired port and any additional settings. This allowed us to use a web browser to view the Jupyter Notebook interface. After launching the Jupyter Notebook interface,

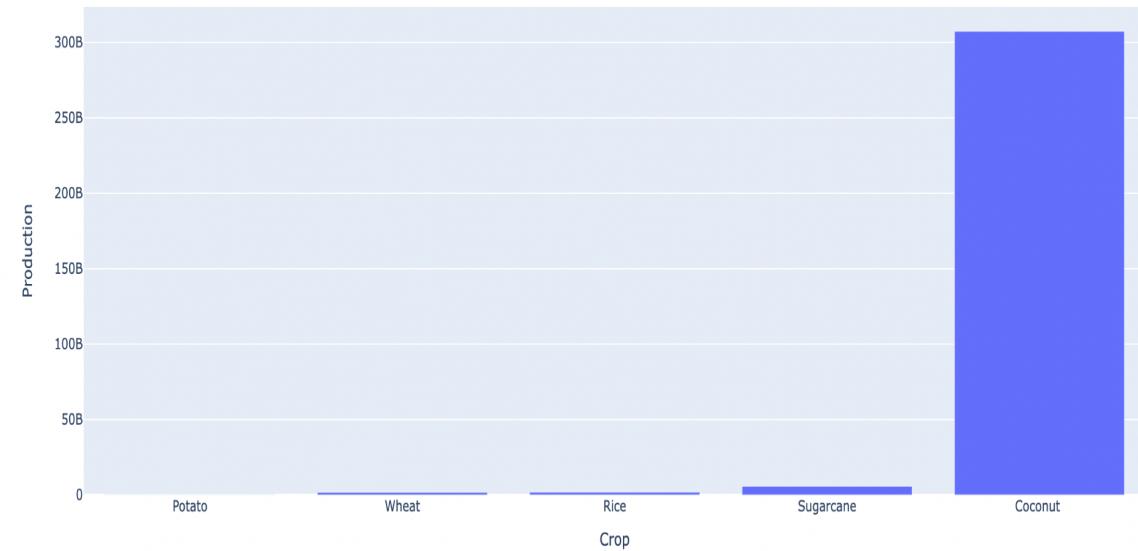
we browsed to the directory holding our.py file containing the data visualization code. We opened the file in Jupyter Notebook, which offered an interactive and user-friendly interface for changing and running the code. With the.py file open, we might use visualization tools such as Matplotlib, Seaborn, or any other to produce meaningful plots, charts, and graphs based on our crop production statistics. We were able to build visual representations of the data by running our data visualization scripts in Jupyter Notebook, allowing us to obtain useful insights and effectively convey our results.

## **GIT LINK:**

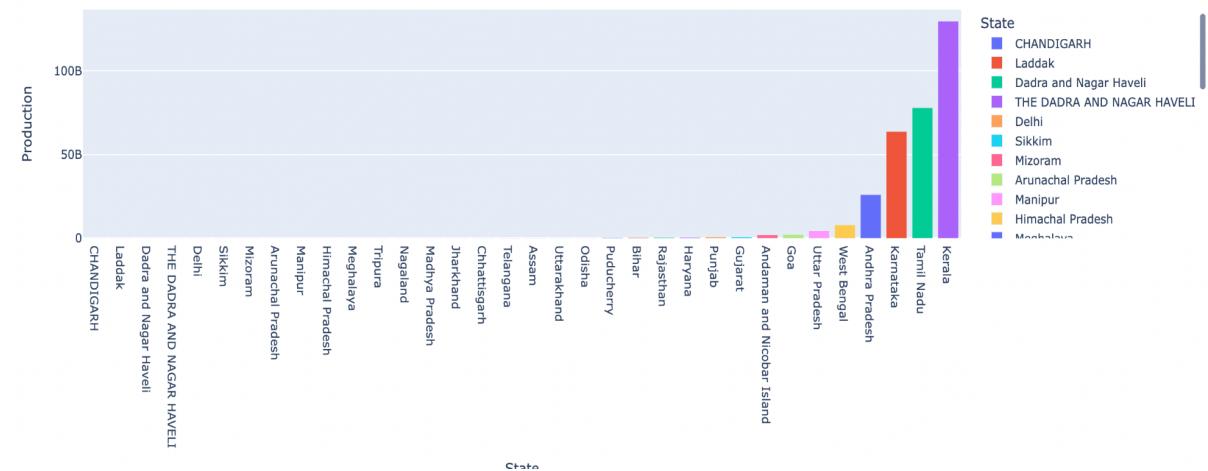
[https://github.com/anuragganjii/DBMS\\_Final\\_Project](https://github.com/anuragganjii/DBMS_Final_Project)

# OUTPUT

```
temp = dataset.groupby(by='Crop')['Production'].sum().reset_index().sort_values(by='Production')
px.bar(temp.tail(), 'Crop', 'Production')
```



```
temp = dataset.groupby(by='State')['Production'].sum().reset_index().sort_values(by='Production')
px.bar(temp, 'State', 'Production', color = 'State')
```

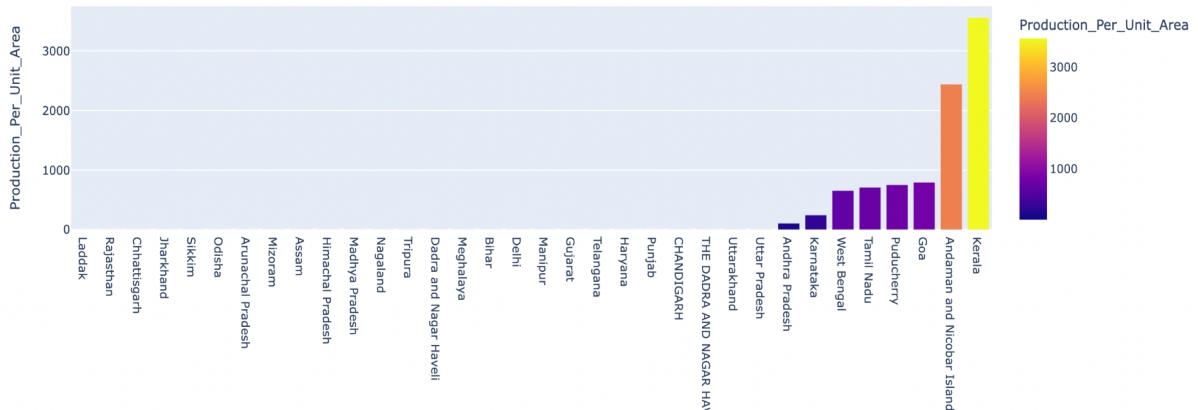


```

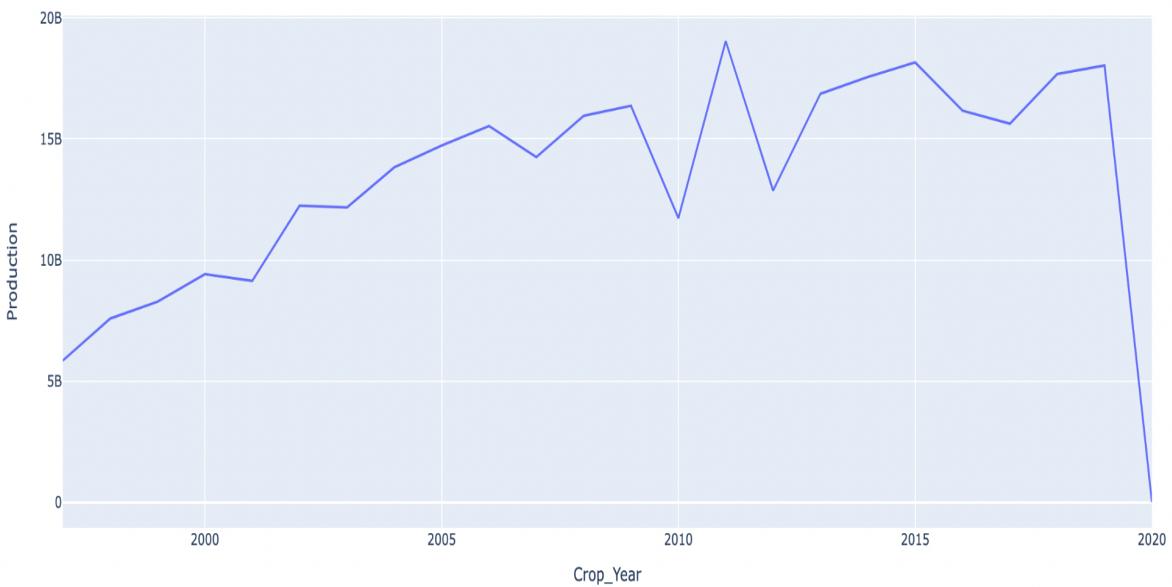
  temp['Production_Per_Unit_Area'] = temp['Production']/temp['Area']
  temp = temp.sort_values(by='Production_Per_Unit_Area')
  px.bar(temp, 'State', 'Production_Per_Unit_Area', color='Production_Per_Unit_Area', )

```

<ipython-input-12-e2bd0ff9c477>:1: FutureWarning:  
Indexing with multiple keys (implicitly converted to a tuple of keys) will be deprecated, use a list instead.

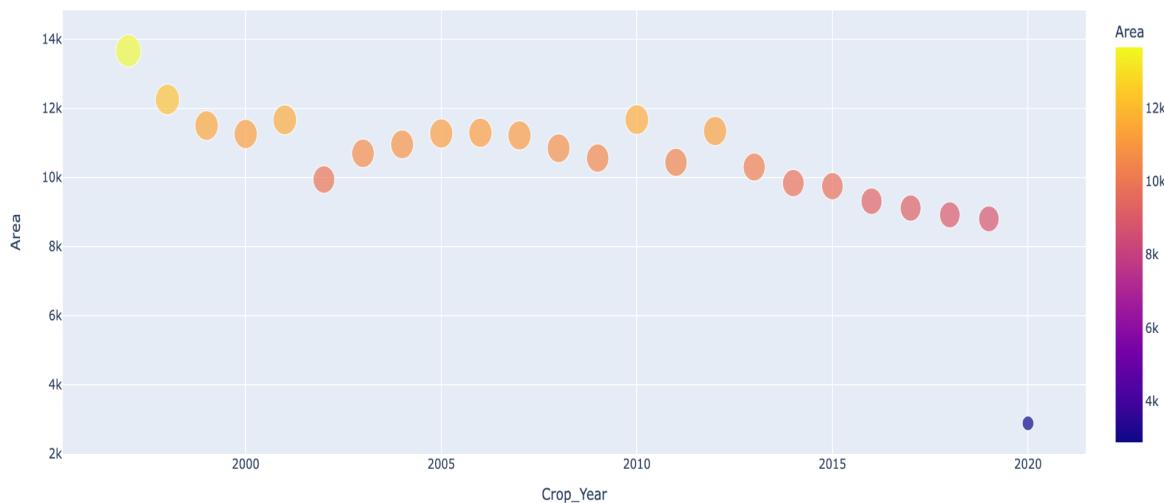


```
[13] temp = dataset.groupby(by='Crop_Year')[['Production']].sum().reset_index()
px.line(temp, 'Crop_Year', 'Production')
```



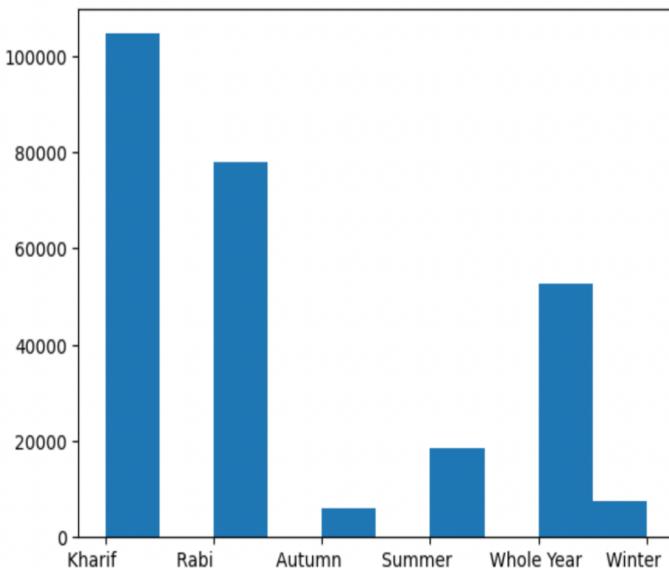
```
temp = dataset.groupby(by='Crop_Year')[['Area']].mean().reset_index()
px.scatter(temp, 'Crop_Year', 'Area', color='Area', size='Area')
```

▶

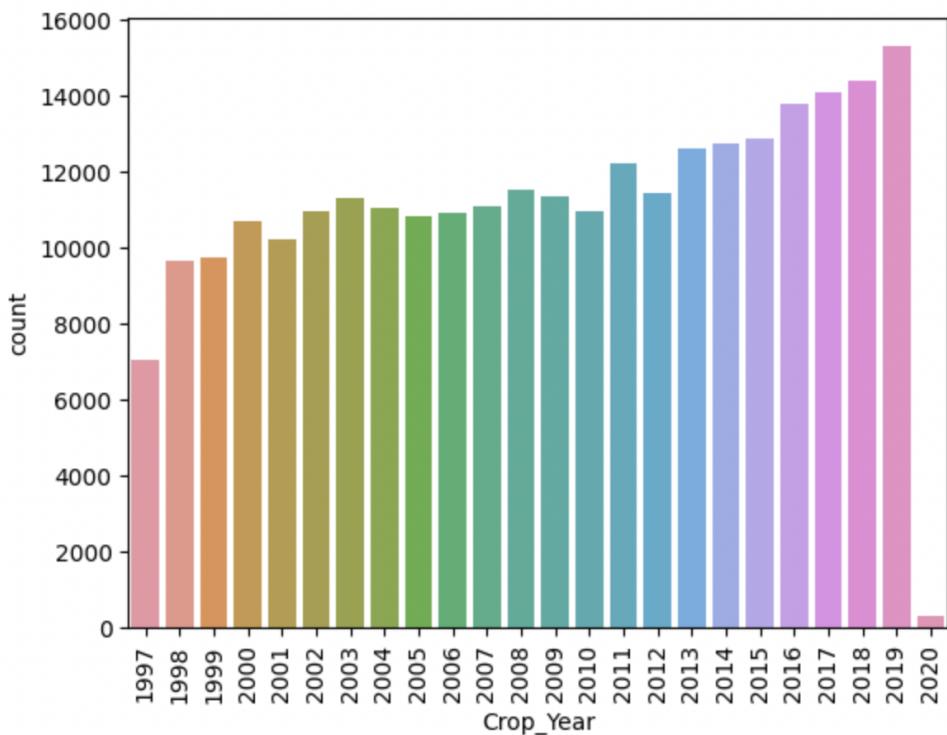


```
plt.hist(dataset['Season'])
```

```
(array([104628., 0., 77980., 0., 6079., 0., 18456.,
       0., 52643., 7583.]),
 array([0., 0.5, 1., 1.5, 2., 2.5, 3., 3.5, 4., 4.5, 5.]),
 <BarContainer object of 10 artists>)
```



```
↳ [Text(21, 0, '2018'),
Text(22, 0, '2019'),
Text(23, 0, '2020'))]
```



# CHALLENGES

We experienced various crop production analysis challenges during our project:

1. Data quality and consistency: We had to assure the dataset's quality and consistency by correcting concerns such as missing data, incorrect entries, and formatting or labeling variances. To ensure the quality and reliability of our study, cleaning and verifying the data needed painstaking attention to detail.
2. Handling vast amounts of data: The agricultural production dataset we worked with was extensive, comprising data from numerous locations, crops, and years. Processing and analyzing such a vast volume of data presented computational and resource problems, requiring us to streamline our processes and utilize technologies such as Spark and BigQuery for effective data processing.
3. Data integration and compatibility: The dataset included data from a variety of sources and formats, necessitating data integration and harmonization. To guarantee compatibility for evaluation, we had to convert, clean, and standardize the data, which needed careful handling and attention to detail.
4. Interpretation and communication: Analyzing crop production data required drawing useful insights and successfully conveying them to agricultural stakeholders and decision-makers. To guarantee that our findings were practical and meaningful, we had to communicate the complicated analytical results in a simple and intelligible manner, using data visualization tools and excellent communication tactics.

# CONCLUSION

Finally, our project focused on predicting agricultural output and gaining insights into crop production trends by employing sophisticated technologies such as Spark, Google Cloud Platform, and machine learning algorithms. We were able to extract useful information from the dataset using data analysis, visualization, and predictive modeling.

Our project's findings have important ramifications for the agriculture business in terms of better planning, resource allocation, and decision-making processes. Farmers may improve their cultivation practices, reduce risks, and boost output by precisely anticipating crop yields. Policymakers may utilize the findings to create successful agriculture policy and promote environmentally friendly practices. Furthermore, the integration of technologies such as Google BigQuery and Tableau Online enabled us to manage, analyze, and visualize the dataset effectively, allowing for better communication of the results. The research highlighted the power of integrating data analysis, machine learning, and data visualization to provide relevant insights and make sound judgments.

Future work might concentrate on broadening the breadth of research, incorporating other data sources, and investigating advanced methodologies to improve forecasts and solve sustainability concerns in agricultural production. Overall, our experiment demonstrated the efficacy of data-driven techniques in agriculture, as well as the possibility for harnessing technology to transform crop production practices. We can pave the road for a more sustainable, efficient, and prosperous future by utilizing the powers of modern data analytics.

## FUTURE WORK

Our study was primarily concerned with predicting crop yield through the application of sophisticated analytics techniques. Moving ahead, various areas of future research may be investigated in order to improve our forecasts and contribute to crop output. To begin, we may look beyond regression into more powerful machine learning models to increase the accuracy of our predictions. Deep learning and ensemble approaches may be used to capture complicated patterns and nonlinear connections in a dataset. Furthermore, including real-time data sources like weather forecasts, satellite imaging, or market data can improve the timeliness and accuracy of our predictions. Creating real-time analytics pipelines and combining streaming data to give up-to-date information for decision-making can be part of this.

Future work might entail expanding the research to include sustainability criteria. We may analyze the environmental effect of crop production and encourage sustainable agricultural practices by taking into account aspects such as water use, pesticide application, and carbon emissions. Overall, our project's future work will include refining our models, including real-time data, examining sustainability aspects, and increasing collaboration within the agricultural sector. By continuing to investigate these areas, we may enhance our forecasts, contribute to crop production technology innovation, and promote sustainable and efficient agricultural practices.