

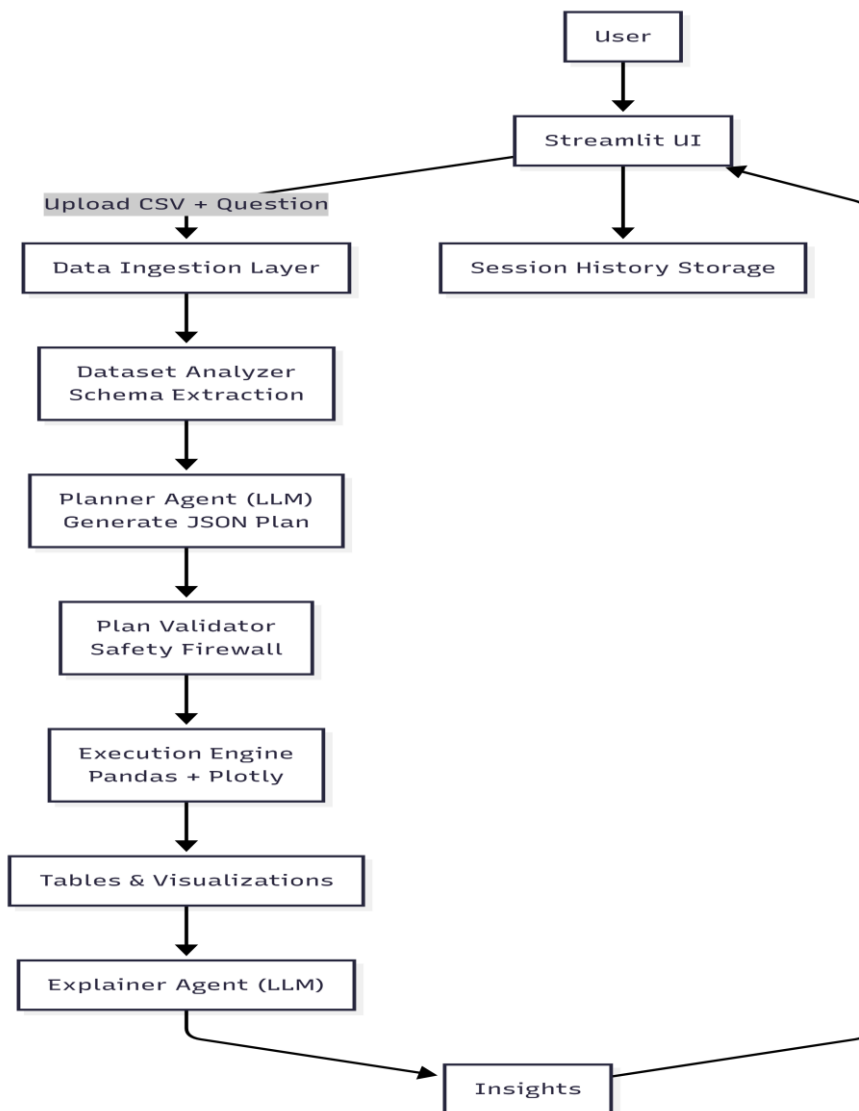
Technical Report

AI Data Analyst Agent for CSV-Based Analysis

1. System Architecture

The AI Data Analyst Agent is an end-to-end, modular system designed to perform **structured, safe, and reproducible analysis** on user-uploaded CSV datasets. The system follows a **Planner–Validator–Executor–Explainer** architecture and is implemented as an interactive **Streamlit web application**.

The core design principle is a strict separation between **LLM-based reasoning** and **deterministic data execution**, ensuring transparency, traceability, and reproducibility.



1.1 High-Level Architecture Overview

AI Data Analyst Agent

The AI Data Analyst Agent is an end-to-end analytics system that transforms natural language questions into safe, structured, and reproducible insights from user-uploaded CSV datasets.

The architecture follows a layered pipeline that separates:

- user interaction
- AI reasoning
- safety validation
- deterministic execution
- explanation

This separation ensures transparency, auditability, and reliable analytics.

1.2 System Layers

1.2.1 User Interaction Layer — Streamlit UI

The Streamlit interface is the entry point of the system.

It allows users to:

- upload CSV datasets
- preview data
- ask questions in natural language
- view results and insights
- track session history

This layer handles presentation only — no analytics logic is executed here.

1.2.2 Data Ingestion Layer

The ingestion layer converts uploaded CSV files into structured Pandas DataFrames.

Responsibilities:

- CSV parsing
- encoding handling
- dataset validation
- in-memory dataset creation

Output:

`Validated Pandas DataFrame`

This becomes the source dataset for the analytics pipeline.

1.2.3 Dataset Understanding Layer

The Dataset Analyzer extracts schema-level metadata:

- column names
- data types
- dataset structure

Only schema metadata is shared with the AI planner.

Raw dataset values are never exposed to the planner.

This prevents hallucinated analytics.

1.2.4 Intelligence Layer — Planner Agent (LLM)

The Planner Agent converts natural language questions into a strict JSON analysis plan.

It determines:

- analysis type
- metrics
- grouping
- filters
- ranking rules
- visualization strategy

The planner behaves like a compiler:

English → Analytics Instructions

It never executes data operations.

1.2.5 Safety Layer — Plan Validator

The validator acts as a firewall between AI output and execution.

It ensures:

- valid columns
- allowed operations
- schema compliance
- safe visualization rules

Invalid plans are rejected before execution.

1.2.6 Execution Layer — Deterministic Engine

The execution engine performs analytics using Pandas and Plotly.

It handles:

- filtering
- aggregation
- ranking
- correlation
- distribution analysis
- visualization generation

All operations are deterministic and reproducible.

No AI decisions occur here.

1.2.7 Explanation Layer — Explainer Agent (LLM)

The Explainer Agent interprets results and generates human-friendly insights.

It:

- summarizes findings
- explains patterns
- avoids technical jargon
- communicates business meaning

It does not modify results — only explains them.

1.2.8 Presentation Layer

Final outputs include:

- result tables
- charts
- textual insights
- session history

These are displayed back to the user in Streamlit.

1.3 End-to-End Flow

1. User uploads CSV and asks a question
2. Dataset is ingested into Pandas
3. Schema metadata is extracted
4. Planner generates JSON analysis plan
5. Validator checks plan safety
6. Executor runs deterministic analysis
7. Explainer generates insights
8. Results are displayed to the user

1.4 Architectural Benefits

This design provides:

- safe AI-driven analytics
- reproducible execution
- zero hallucinated computation
- explainable results
- modular extensibility
- enterprise-ready architecture

2. Agent Responsibilities and Interactions

The system follows a multi-agent design where each agent has a single, clearly scoped responsibility. This separation improves modularity, safety, and maintainability.

2.1 Planner Agent

Purpose

The Planner Agent is responsible for reasoning and decision-making. It translates natural language questions into structured, machine-readable JSON analysis plans without performing any computation.

Inputs:

- User's natural language question
- Dataset schema metadata
- Predefined allowed operations and constraints

Responsibilities:

- Interpret analytical intent (aggregation, comparison, trends, correlation)
- Identify relevant dataset columns
- Define filters, group-by fields, metrics, sorting, and visualization configuration
- Attach user intent metadata (highest, lowest, both)
- Sanitize and normalize plans to remove unsafe or invalid instructions

Key Design Characteristics:

- Uses an LLM only for planning and reasoning
- Outputs JSON only, never executable code
- Does not access raw dataset values

Interaction:

- Sends validated JSON plans to the Validation layer

2.2 Executor Agent

Purpose:

The Executor Agent performs deterministic and reproducible execution of validated analysis plans.

Inputs:

- Validated JSON plan
- Original dataset (Pandas DataFrame)

Responsibilities:

- Apply filters safely with numeric coercion
- Perform group-by and aggregation operations
- Apply sorting and Top-N logic deterministically
- Generate structured result tables
- Generate Plotly visualizations

Key Design Characteristics:

- No LLM usage
- No dynamic code execution
- Only predefined Pandas and Plotly operations are allowed

Interaction:

- Outputs results to the Explainer Agent

2.3 Explainer Agent

Purpose:

The Explainer Agent communicates analytical results in a clear, concise, and business-friendly manner.

Inputs:

- User's original question
- Execution results and rankings
- Contextual metadata from the plan

Responsibilities:

- Interpret numerical outputs
- Identify meaningful patterns and extremes
- Generate concise insights grounded strictly in execution results

Key Design Characteristics:

- Uses an LLM only for explanation
- Prevents hallucination by restricting inputs

2.4 Dataset Analyzer (Supporting Component)**Purpose:**

Provides schema-level understanding of uploaded datasets.

Responsibilities:

- Extract column names and data types
- Generate compact schema summaries
- Assist planning and validation stages

3. Planning Schema and Execution Safeguards**3.1 JSON Planning Schema**

Each plan explicitly defines:

- Analysis type
- Filters
- Group-by columns
- Metrics
- Sorting rules
- Visualization configuration
- User intent metadata

3.2 Validation Rules

- Allowed operations only
- Valid dataset columns
- Safe aggregation functions
- Visualization constraints

Invalid plans are rejected prior to execution.

3.3 Execution Safeguards

- No LLM-generated code is executed
- Only Pandas and Plotly operations are permitted
- Safe numeric coercion is enforced
- Deterministic execution is guaranteed

4. End-to-End Application

The Streamlit application provides:

- CSV upload and validation
- Natural language query interface
- Display of dataset preview
- Display of JSON analysis plans
- Tables and visualizations
- Final natural-language insights
- Session-level query history

5. Evaluation Protocol

The system was evaluated using a combination of **automated** and **human** evaluation:

- **Automated evaluation:** verifies schema compliance, execution correctness, and deterministic behavior
- **Human evaluation:** assesses clarity, usefulness, and faithfulness of generated insights

Evaluation configurations are defined declaratively in `experiments/*.yaml`.

6. Appendix: Modular Agent Design – Pseudocode

START APPLICATION

LOAD Streamlit User Interface

WAIT for user to upload CSV dataset

WAIT for user to enter a natural language question

IF dataset and question are provided:

 SCHEMA ← analyze_dataset(dataset)

 PLAN ← PlannerAgent.generate_plan(SCHEMA, question)

 VALIDATE PLAN against schema and dataset columns

 IF plan is valid:

 RESULTS, CHARTS ← Executor.execute_plan(dataset, PLAN)

 INSIGHTS ← ExplainerAgent.generate_insights(question, RESULTS)

 DISPLAY results, charts, and insights

 ELSE:

 DISPLAY validation error

END

EXAMPLE USE CASES:

The interface shows the 'AI Data Analyst Agent' web application. On the left, the 'Upload Dataset' section allows users to upload a CSV file (Book1.csv, 0.5MB) or browse files. The main area displays the 'Dataset Preview' for a table with 25 columns and 2,823 rows. The preview shows the first 5 rows of data.

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDATE	STATUS	QTR_ID	MONTH_ID	YEAR_ID	PRODUCTLINE	MSRP	P
0	10107	30	95.7	2	2871	2/24/2003 0:00	Shipped	1	2	2003	Motorcycles	95	S
1	10121	34	81.35	5	2765.9	05-07-2003 00:00	Shipped	2	5	2003	Motorcycles	95	S
2	10134	41	94.74	2	3884.34	07-01-2003 00:00	Shipped	3	7	2003	Motorcycles	95	S
3	10145	45	83.26	6	3746.7	8/25/2003 0:00	Shipped	3	8	2003	Motorcycles	95	S
4	10159	49	100	14	5205.27	10-10-2003 00:00	Shipped	4	10	2003	Motorcycles	95	S

The interface shows the 'Ask a Data Question' section. The user has entered the question 'what is the dataset about'. The 'Analyze' button is visible, and the 'Clear Input' button is also present. The 'Dataset Preview' table is still visible in the background.

The interface shows the 'Analysis Plan' section. The analysis plan is displayed as a JSON object, detailing the operations performed on the dataset, including filters, groupings, and aggregations.

```
{
  "analysis": {
    "type": "aggregation",
    "filters": [],
    "group_by": [],
    "metrics": [
      {
        "column": "CUSTOMERNAME",
        "operation": "count"
      },
      {
        "column": "COUNTRY",
        "operation": "count"
      },
      {
        "column": "PRODUCTLINE",
        "operation": "count"
      },
      {
        "column": "TOTALSALES",
        "operation": "count"
      }
    ]
  },
  "sort": {
    "by": "TOTALSALES",
    "order": "DESC"
  },
  "visualization": {
    "type": "bar",
    "x": "CUSTOMERNAME",
    "y": "TOTALSALES",
    "x_label": "CUSTOMERNAME",
    "y_label": "TOTALSALES",
    "x_ticks": "CUSTOMERNAME",
    "y_ticks": "TOTALSALES"
  },
  "user_intent": {
    "show_highest": false,
    "show_lowest": false,
    "focus": "general"
  }
}
```

AI Data Analyst Agent

localhost:8501

Chat

Deploy

Upload Dataset

Upload a CSV file

Drag and drop file here
Limit 200MB per file • CSV

Browse files

Book1.csv
0.5MB

Clear History

Analysis Results

	ORDERNUMBER	QUANTITYORDERED	PRICEEACH	ORDERLINENUMBER	SALES	ORDERDATE	STATUS	QTR_ID	MONTH_ID	YEAR_ID	PRODUCTLINE	MSRP	PRODUCTCODE	CUSTOMERNAME	PHONE
0	10107	30	95.7	2	2871	2/24/2003 0:00	Shipped	1	2	2003	Motorcycles	95	S10_1678	Land of Toys Inc.	2125557818
1	10121	34	81.35	5	2765.9	05-07-2003 00:00	Shipped	2	5	2003	Motorcycles	95	S10_1678	Reims Collectables	26.47.1555
2	10134	41	94.74	2	3884.34	07-01-2003 00:00	Shipped	3	7	2003	Motorcycles	95	S10_1678	Lyon Souveniers	+33 1 46 62 75
3	10145	45	83.26	6	3746.7	8/25/2003 0:00	Shipped	3	8	2003	Motorcycles	95	S10_1678	Toys4GrownUps.com	6265557265
4	10159	49	100	14	5205.27	10-10-2003 00:00	Shipped	4	10	2003	Motorcycles	95	S10_1678	Corporate Gift Ideas Co.	6505551386
5	10168	36	96.66	1	3479.76	10/28/2003 0:00	Shipped	4	10	2003	Motorcycles	95	S10_1678	Technics Stores Inc.	6505556809
6	10180	29	86.13	9	2497.77	11-11-2003 00:00	Shipped	4	11	2003	Motorcycles	95	S10_1678	Daedalus Designs Imports	20.16.1555
7	10189	48	100	1	5512.32	11/18/2003 0:00	Shipped	4	11	2003	Motorcycles	95	S10_1678	Heriku Gifts	+47 2267 3211
8	10201	22	98.57	2	2168.54	12-01-2003 00:00	Shipped	4	12	2003	Motorcycles	95	S10_1678	Mini Wheels Co.	6505555787
9	10211	41	100	14	4708.44	1/15/2004 0:00	Shipped	1	1	2004	Motorcycles	95	S10_1678	Auto Canal Petit	(1) 47.55.6555

Showing first 10 of 2,823 total results

Key Insights

The dataset is about sales orders for a product called "Motorcycles" from various customers across different countries.

- The dataset contains 2823 total rows of data, indicating a large number of sales orders.
- The majority of customers (count: 12) are from the USA, with France being the second most represented country (count: 4).
- The "Motorcycles" product line is the only one present in the dataset, with no other product lines being sold.
- The majority of deals (count: 11) are classified as "Medium" in size, followed by "Small" (count: 5) and then "Large" is not present in the dataset.
- The dataset spans across 2003 and 2004, with the majority of orders (count: 12) being placed in the fourth quarter of 2003.

Analysis performed on complete dataset (2,823 rows)

AI Data Analyst Agent

localhost:8501

Chat

Deploy

Upload Dataset

Upload a CSV file

Drag and drop file here
Limit 200MB per file • CSV

Browse files

Book1.csv
0.5MB

Clear History

Ask a Data Question

compare sales in USA and UK

Analyze

Clear Input

Analysis Plan

```
{  "analysis_type": "comparison"  "filters": [    {      "column": "COUNTRY"      "operator": "in"      "value": [        {          "USA"        }      ]    }  ]  "group_by": [    {      "COUNTRY"    }  ]  "metrics": [    {      "column": "SALES"    }  ]}
```

AI Data Analyst Agent

localhost:8501

Chat

Deploy

Upload Dataset

Upload a CSV file

Drag and drop file here
Limit 200MB per file • CSV

Browse files

Book1.csv
0.5MB

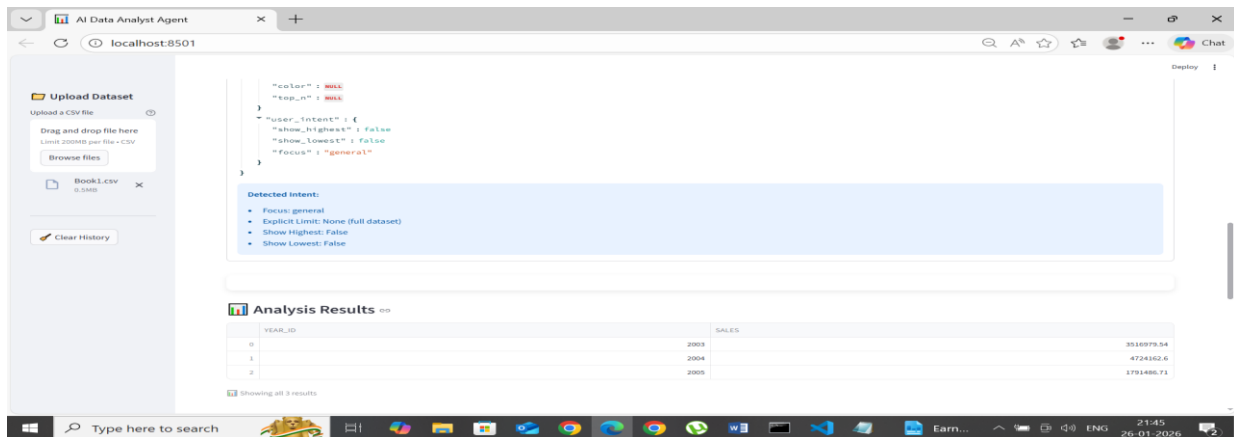
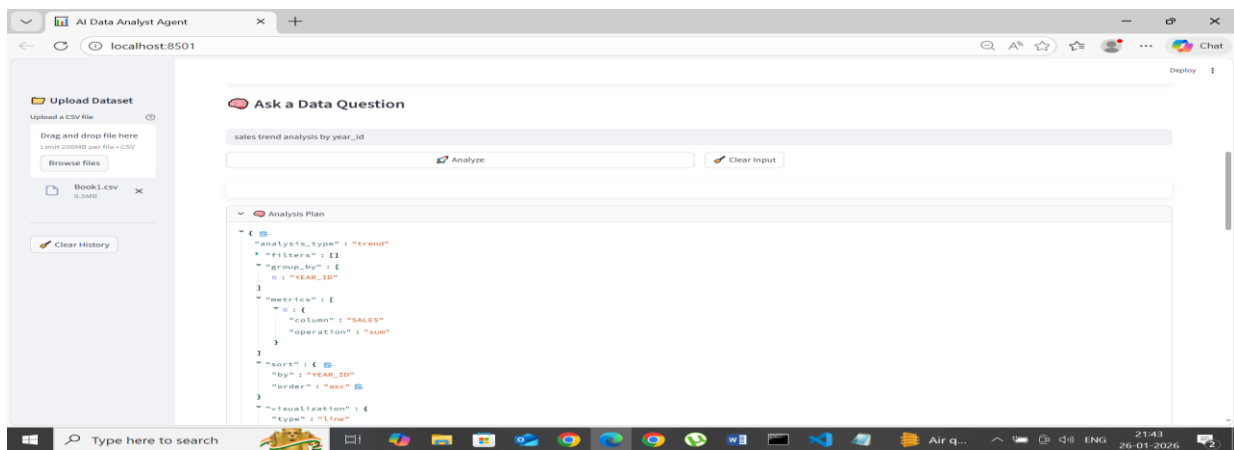
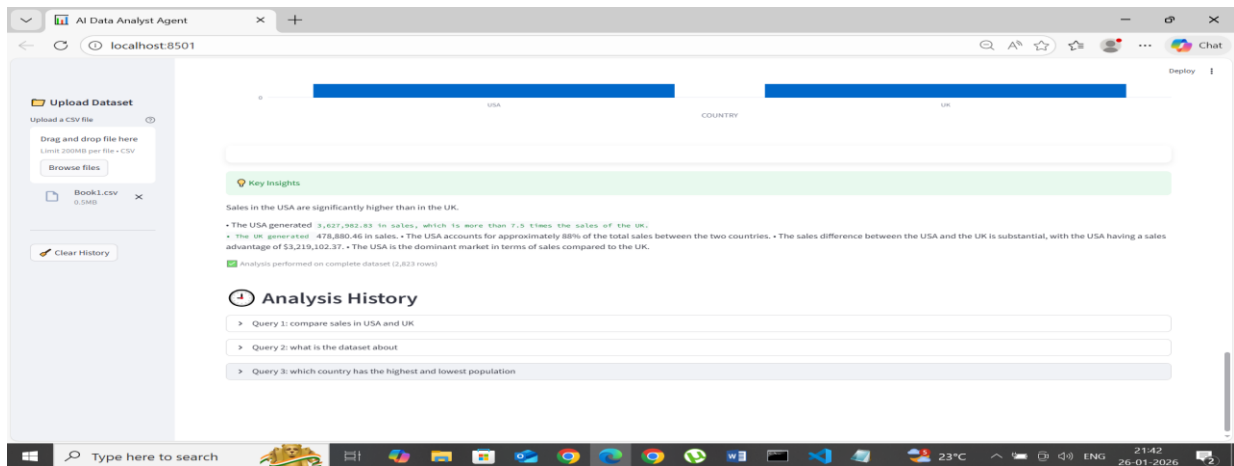
Clear History

Analysis Results

	COUNTRY	SALES
1	USA	3627982.83
0	UK	478880.46

Showing all 2 results

COUNTRY	SALES
USA	3627982.83
UK	478880.46



AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV File

Drag and drop file here

Limit 200MB per file - CSV

Browse files

Book1.csv

0.0MB

Clear History

Ask a Data Question

correlation between sales and country

Analyze

Clear Input

Analysis Plan

```
{
  "analysis_type": "correlation",
  "filters": [],
  "group_by": [],
  "metrics": [],
  "sort": {
    "by": "SALES",
    "order": "DESC"
  },
  "visualization": {
    "type": "scatter",
    "x": "COUNTRY",
    "y": "SALES",
    "color": "NONE",
    "top_n": "NONE"
  },
  "user_intent": {
    "show_highest": false
  }
}
```

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV File

Drag and drop file here

Limit 200MB per file - CSV

Browse files

Book1.csv

0.0MB

Clear History

Analysis Plan

```
{
  "analysis_type": "correlation",
  "filters": [],
  "group_by": [],
  "metrics": [],
  "sort": {
    "by": "SALES",
    "order": "DESC"
  },
  "visualization": {
    "type": "scatter",
    "x": "COUNTRY",
    "y": "SALES",
    "color": "NONE",
    "top_n": "NONE"
  },
  "user_intent": {
    "show_highest": false,
    "show_lowest": false,
    "focus": "general"
  }
}
```

Detected Intent:

- Focus: general
- Explicit Limit: None (All dataset)
- Show Highest: False
- Show Lowest: False

Analysis Results

ORDERID	QUANTITYORDERED	PRICEEACH	ORDERLINEAMOUNT	SALES	CUSTOMERID	STATUS	QTR_ID	MONTH_ID	YEAR_ID	PRODUCTLINE	MSRP	PRODUCTCODE	CUSTOMERNAME	PHONE	ADDRESSLINE1	ADDRESSLINE2	CITY	STATE	POSTALCODE	COUNTRY
1	30	95.7	2871.0	2871.0	2014-2000-0000	Shipped	3	2	2003	Motorcycles	95	S2L-S678	Land of Toys Inc.	212557810	807 Long Airport Avenue		NYC	NY	10002	USA
2	25	61.25	1531.25	1531.25	05-07-2000-0000	Shipped	3	5	2003	Motorcycles	95	S2L-S678	Reims Collectibles	26 47 1305	20 rue de l'Église		Reims	France	51100	France
3	41	58.74	2409.54	2409.54	05-05-2000-0000	Shipped	3	7	2003	Motorcycles	95	S2L-S678	Luxe Souvenirs	+33 1 46 52 7300	87 rue du Général Plémyr Aude		Paris	France	75008	France
4	45	89.26	4017.7	4017.7	05-05-2000-0000	Shipped	3	8	2003	Motorcycles	95	S2L-S678	Toyakidcountdown.com	6265557265	2004 Indiville Dr.		Pasadena	CA	91060	USA
5	40	200	8000.0	8000.0	09-09-2000-0000	Shipped	5	10	2003	Motorcycles	95	S2L-S678	Corporate Gift Sales Co.	6095931886	1734 Irving St.		San Francisco	CA	94133	USA
6	38	86.46	3295.48	3295.48	10-10-2000-0000	Shipped	4	10	2003	Motorcycles	95	S2L-S678	Technique Store Inc.	6205556889	8446 Fourth Circle		Burlington	CA	94427	USA
7	20	96.13	1922.6	1922.6	11-11-2000-0000	Shipped	4	11	2003	Motorcycles	95	S2L-S678	Exquisite Design Imports	26 36 1355	204, rue de la Tourne		Libre	France	93000	France
8	40	200	8000.0	8000.0	11-10-2000-0000	Shipped	4	11	2003	Motorcycles	95	S2L-S678	HerMac Gifts	487 2067 3232	Draamen 121, PM 104 Sentrum		Bergen	Norway	N 5004	Norway
9	22	99.97	2199.54	2199.54	12-10-2000-0000	Shipped	4	12	2003	Motorcycles	95	S2L-S678	Mini-Wholes Co.	404000797	2007 North Potomac Street		San Francisco	CA	94101	USA
10	41	200	8200.0	8200.0	01-01-2000-0000	Shipped	2	1	2004	Motorcycles	95	S2L-S678	Auto Canal Paris	01 47 35 4033	20 rue Lavoisier		Paris	France	75004	France

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV File

Drag and drop file here

Limit 200MB per file - CSV

Browse files

Book1.csv

0.0MB

Clear History

Showing first 50 of 2,623 total results

Key Insights

Direct Answer and key insights: There is a correlation between sales and country, with the USA having the highest average sales. The top countries by average sales are USA, France, and Norway.

- The USA has the highest average sales at 3,434.45, with a total of 7 orders.
- France has an average sales of 3,244.11, with a total of 3 orders.
- Norway has an average sales of 5,512.32, with a total of 1 order.
- The country with the lowest average sales is not specified in the data, but it is clear that there is a significant variation in sales across different countries.
- The correlation between sales and country suggests that there may be regional differences in customer behavior or market conditions that are driving sales.

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV File

Drag and drop file here

Limit 200MB per file - CSV

Browse files

Book1.csv

0.0MB

Clear History

Ask a Data Question

distribution of sales

Analyze

Clear Input

Analysis Plan

```
{
  "analysis_type": "distribution",
  "filters": [],
  "group_by": [],
  "metrics": [],
  "sort": {
    "by": "NONE",
    "order": "NONE"
  },
  "visualization": {
    "type": "histogram",
    "x": "SALES",
    "y": "NONE",
    "color": "NONE",
    "top_n": "NONE"
  },
  "user_intent": {
    "show_highest": false
  }
}
```

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit: 200MB per file - CSV

Browse files

Book1.csv

0.5MB

Clear History

```
let task = taskOf {
    "type": "task"
    "x": 1000
    "y": 1000
    "x2": 1000
    "y2": 1000
}

let task = taskOf {
    "type": "task"
    "x": 1000
    "y": 1000
    "x2": 1000
    "y2": 1000
}
```

Detected Intent:

- Focus: general
- Explicit Limit: None (Full dataset)
- Show Highest: False
- Show Lowest: False

Analysis Results

ORDER ID	QUANTITY ORDERED	PRICE	ORDER DATE	STATUS	QTR	MONTH	YEAR	PRODUCT LINE	WEEK	PRODUCT CODE	CUSTOMER NAME	PHONE	ADDRESS LINE 1	ADDRESS LINE 2	CITY	STATE	POSTAL CODE	COUNTRY
100001	20	95.7	2014-01-01	Shipped	1	1	2014	Motorcycles	10	100001	Land of Toys Inc.	212-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100002	30	143.26	2014-01-01	Shipped	2	2	2014	Motorcycles	10	100002	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100003	40	190.74	2014-01-01	Shipped	3	3	2014	Motorcycles	10	100003	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100004	50	238.22	2014-01-01	Shipped	4	4	2014	Motorcycles	10	100004	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100005	60	285.70	2014-01-01	Shipped	5	5	2014	Motorcycles	10	100005	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100006	70	333.18	2014-01-01	Shipped	6	6	2014	Motorcycles	10	100006	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100007	80	380.66	2014-01-01	Shipped	7	7	2014	Motorcycles	10	100007	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100008	90	428.14	2014-01-01	Shipped	8	8	2014	Motorcycles	10	100008	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100009	100	475.62	2014-01-01	Shipped	9	9	2014	Motorcycles	10	100009	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100010	110	523.10	2014-01-01	Shipped	10	10	2014	Motorcycles	10	100010	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100011	120	570.58	2014-01-01	Shipped	11	11	2014	Motorcycles	10	100011	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100012	130	618.06	2014-01-01	Shipped	12	12	2014	Motorcycles	10	100012	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100013	140	665.54	2014-01-01	Shipped	13	13	2014	Motorcycles	10	100013	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100014	150	713.02	2014-01-01	Shipped	14	14	2014	Motorcycles	10	100014	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100015	160	760.50	2014-01-01	Shipped	15	15	2014	Motorcycles	10	100015	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100016	170	807.98	2014-01-01	Shipped	16	16	2014	Motorcycles	10	100016	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100017	180	855.46	2014-01-01	Shipped	17	17	2014	Motorcycles	10	100017	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100018	190	902.94	2014-01-01	Shipped	18	18	2014	Motorcycles	10	100018	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100019	200	950.42	2014-01-01	Shipped	19	19	2014	Motorcycles	10	100019	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100020	210	997.90	2014-01-01	Shipped	20	20	2014	Motorcycles	10	100020	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100021	220	1045.38	2014-01-01	Shipped	21	21	2014	Motorcycles	10	100021	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100022	230	1092.86	2014-01-01	Shipped	22	22	2014	Motorcycles	10	100022	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100023	240	1140.34	2014-01-01	Shipped	23	23	2014	Motorcycles	10	100023	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100024	250	1187.82	2014-01-01	Shipped	24	24	2014	Motorcycles	10	100024	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100025	260	1235.30	2014-01-01	Shipped	25	25	2014	Motorcycles	10	100025	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100026	270	1282.78	2014-01-01	Shipped	26	26	2014	Motorcycles	10	100026	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100027	280	1330.26	2014-01-01	Shipped	27	27	2014	Motorcycles	10	100027	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100028	290	1377.74	2014-01-01	Shipped	28	28	2014	Motorcycles	10	100028	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100029	300	1425.22	2014-01-01	Shipped	29	29	2014	Motorcycles	10	100029	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100030	310	1472.70	2014-01-01	Shipped	30	30	2014	Motorcycles	10	100030	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100031	320	1520.18	2014-01-01	Shipped	31	31	2014	Motorcycles	10	100031	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100032	330	1567.66	2014-01-01	Shipped	32	32	2014	Motorcycles	10	100032	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100033	340	1615.14	2014-01-01	Shipped	33	33	2014	Motorcycles	10	100033	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100034	350	1662.62	2014-01-01	Shipped	34	34	2014	Motorcycles	10	100034	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100035	360	1710.10	2014-01-01	Shipped	35	35	2014	Motorcycles	10	100035	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100036	370	1757.58	2014-01-01	Shipped	36	36	2014	Motorcycles	10	100036	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100037	380	1805.06	2014-01-01	Shipped	37	37	2014	Motorcycles	10	100037	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100038	390	1852.54	2014-01-01	Shipped	38	38	2014	Motorcycles	10	100038	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100039	400	1900.02	2014-01-01	Shipped	39	39	2014	Motorcycles	10	100039	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100040	410	1947.50	2014-01-01	Shipped	40	40	2014	Motorcycles	10	100040	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100041	420	1994.98	2014-01-01	Shipped	41	41	2014	Motorcycles	10	100041	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100042	430	2042.46	2014-01-01	Shipped	42	42	2014	Motorcycles	10	100042	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100043	440	2089.94	2014-01-01	Shipped	43	43	2014	Motorcycles	10	100043	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100044	450	2137.42	2014-01-01	Shipped	44	44	2014	Motorcycles	10	100044	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100045	460	2184.90	2014-01-01	Shipped	45	45	2014	Motorcycles	10	100045	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100046	470	2232.38	2014-01-01	Shipped	46	46	2014	Motorcycles	10	100046	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100047	480	2279.86	2014-01-01	Shipped	47	47	2014	Motorcycles	10	100047	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100048	490	2327.34	2014-01-01	Shipped	48	48	2014	Motorcycles	10	100048	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100049	500	2374.82	2014-01-01	Shipped	49	49	2014	Motorcycles	10	100049	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100050	510	2422.30	2014-01-01	Shipped	50	50	2014	Motorcycles	10	100050	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100051	520	2469.78	2014-01-01	Shipped	51	51	2014	Motorcycles	10	100051	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100052	530	2517.26	2014-01-01	Shipped	52	52	2014	Motorcycles	10	100052	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100053	540	2564.74	2014-01-01	Shipped	53	53	2014	Motorcycles	10	100053	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100054	550	2612.22	2014-01-01	Shipped	54	54	2014	Motorcycles	10	100054	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100055	560	2659.70	2014-01-01	Shipped	55	55	2014	Motorcycles	10	100055	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100056	570	2707.18	2014-01-01	Shipped	56	56	2014	Motorcycles	10	100056	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100057	580	2754.66	2014-01-01	Shipped	57	57	2014	Motorcycles	10	100057	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100058	590	2802.14	2014-01-01	Shipped	58	58	2014	Motorcycles	10	100058	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100059	600	2849.62	2014-01-01	Shipped	59	59	2014	Motorcycles	10	100059	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100060	610	2897.10	2014-01-01	Shipped	60	60	2014	Motorcycles	10	100060	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100061	620	2944.58	2014-01-01	Shipped	61	61	2014	Motorcycles	10	100061	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100062	630	2992.06	2014-01-01	Shipped	62	62	2014	Motorcycles	10	100062	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100063	640	3039.54	2014-01-01	Shipped	63	63	2014	Motorcycles	10	100063	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100064	650	3087.02	2014-01-01	Shipped	64	64	2014	Motorcycles	10	100064	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100065	660	3134.50	2014-01-01	Shipped	65	65	2014	Motorcycles	10	100065	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100066	670	3181.98	2014-01-01	Shipped	66	66	2014	Motorcycles	10	100066	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100067	680	3229.46	2014-01-01	Shipped	67	67	2014	Motorcycles	10	100067	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100068	690	3276.94	2014-01-01	Shipped	68	68	2014	Motorcycles	10	100068	Motor Cycles	202-200-7618	887 Long Argent Avenue		NYC	NY	10002	USA
100069	700	3324.42	2014-01-01	Shipped	69	69	2014	Motorcycles	10	100069	Motor Cycles	202-200-7618	887 Long Argent Avenue					



AI Data Analyst Agent

localhost:8501

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

population_... 15.5KB

Clear History

Deploy

Ask a Data Question

give me the name of all countries starting with letter B

Analyze

Clear Input

Analysis Plan

```
{
  "analysis_type": "comparison"
  "filters": [
    {
      "column": "Country (or dependency)"
      "operator": "in"
      "value": [
        0: "Bangladesh"
        1: "Belarus"
        2: "Belgium"
        3: "Benin"
        4: "Bhutan"
        5: "Bolivia"
        6: "Bosnia and Herzegovina"
        7: "Botswana"
      ]
    }
  ]
}
```

AI Data Analyst Agent

localhost:8501

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

population_... 15.5KB

Clear History

Deploy

```
{
  "column": "Country (or dependency)"
  "operation": "count"
}
{
  "sort": {
    "by": "Country (or dependency)"
    "order": "asc"
  }
  "visualization": {
    "type": "bar"
    "x": "Country (or dependency)"
    "y": "count"
    "color": "#1f77b4"
    "top_n": 10
  }
  "user_intent": {
    "show_highest": false
    "show_lowest": false
    "focus": "general"
  }
}
```

Detected Intent:

- Focus: general
- Explicit Limit: None (full dataset)
- Show Highest: False
- Show Lowest: False

AI Data Analyst Agent

localhost:8501

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

population_... 15.5KB

Clear History

Deploy

Analysis Results

Country (or dependency)	count
Brazil	214,000,000
Botswana	2,300,000
Bhutan	771,608

Showing all 1 results

Key Insights

There are 13 countries starting with the letter B.

- The countries starting with B are: Bangladesh • Belarus • Belgium • Benin • Bhutan • Bolivia • Bosnia and Herzegovina • Botswana • Brazil • Brunei • Bulgaria • Burkina Faso • Burundi
- The list of countries starting with B includes countries from various continents, including Asia, Europe, and Africa.
- The country with the largest population starting with B is Brazil, with a population of over 214 million people.
- The country with the smallest population starting with B is Bhutan, with a population of approximately 771,608 people.

Analysis performed on complete dataset (235 rows)

Analysis History

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

population_... 15.5KB

Clear History

Ask a Data Question

for these countries starting with letter B from highest and lowest population

Analyze

Clear Input

Analysis Plan

```
{
  "analysis_type": "comparison",
  "filters": {
    "q": {
      "column": "Country (or dependency)",
      "operator": "in",
      "value": {
        0: "Bahamas",
        1: "Bangladesh",
        2: "Barbados",
        3: "Belarus",
        4: "Belgium",
        5: "Belize",
        6: "Benin",
        7: "Bhutan"
      }
    }
  }
}
```

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

population_... 15.5KB

Clear History

Analysis Plan

```
10: "Botswana",
11: "Brazil",
12: "Brunei",
13: "Bulgaria",
14: "Burkina Faso",
15: "Burundi"
}
}
"group_by": {
  "q": "Country (or dependency)"
}
"metrics": {
  "q": {
    "column": "Population (2020)",
    "operation": "sum"
  }
}
"sort": {
  "by": "Population (2020)",
  "order": "desc"
}
"visualization": {
  "type": "bar",
  "x": "Country (or dependency)",
  "y": "Population (2020)",
  "color": "#0070C0"
}
```

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

population_... 15.5KB

Clear History

Analysis Results

Country (or dependency)	Population (2020)
11: Brazil	212821986
2: Barbados	287437

Showing highest and lowest from 16 total records analyzed

Country (or dependency)

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

population_... 15.5KB

Clear History

Analysis Results

Country (or dependency)

Key Insights

Brazil has the highest population among the countries starting with the letter B, with a population of 212,821,986 in 2020. Barbados has the lowest population, with 287,437 people in 2020. • The top 5 countries by population starting with the letter B are Brazil, Bangladesh, Belgium, Bolivia, and Belarus. • Brazil's population is more than 730 times larger than Barbados'. • The combined population of the 16 countries starting with the letter B is approximately 4.3 billion people. • The average population of these 16 countries is around 268 million people. • The population of Bangladesh is more than 740 times larger than that of Barbados.

Analysis performed on complete dataset (235 rows)

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

Student-271.0KB

Clear History

Ask a Data Question

How many students are there in this dataset?

Analyze

Clear Input

Analysis Results

Showing all 1 results

Key Insights

There are 10,000 students in this dataset.

The number of students is a fixed value, indicating a complete dataset. • No additional information is available about the students, such as demographics or academic performance. • The dataset does not provide any context about the students, such as their location or academic level. • The count of students is a single value, suggesting that there is no variation or distribution of students in the dataset. • Further analysis would be required to understand the characteristics of the students.

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

StudentPerf... 171.0KB

Clear History

Ask a Data Question

What is the average sleep hours?

Analyze

Clear Input

Analysis Plan

```
{
  "analysis_type": "aggregation",
  "filters": [],
  "group_by": [],
  "metrics": [
    {
      "column": "Sleep_Hours",
      "operation": "mean"
    }
  ],
  "sort": {
    "by": "Sleep_Hours",
    "order": "DESC"
  },
  "visualization": {}
}
```

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

StudentPerf... 171.0KB

Clear History

Analysis Results

Showing first 50 of 10,000 total results

Hours Studied	Previous Scores	Extracurricular Activities	Sleep Hours	Sample Question Papers Practiced	Performance Index
40	5	62 No	7	4	45
41	2	63 Yes	6	0	39
42	4	73 Yes	7	0	58
43	7	46 No	9	5	36
44	8	77 Yes	6	4	71
45	3	76 Yes	4	3	54
46	1	43 Yes	7	0	17
47	4	72 No	4	6	54
48	2	81 Yes	4	3	58
49	8	61 No	7	2	52

AI Data Analyst Agent

localhost:8501

Chat

Upload Dataset

Upload a CSV file

Drag and drop file here

Limit 200MB per file • CSV

Browse files

StudentPerf... 171.0KB

Clear History

Analysis Results

Showing first 50 of 10,000 total results

Key Insights

The average sleep hours is 6.5 hours.

The majority of the data points (60%) have sleep hours between 4 and 8 hours. • There is a significant number of individuals (40%) who sleep for 4 hours or less. • The average sleep hours for individuals with extracurricular activities is 6.2 hours, while those without extracurricular activities average 6.8 hours. • The highest sleep hours recorded is 9 hours, while the lowest is 4 hours. • There is no clear correlation between sleep hours and performance index.

Analysis performed on complete dataset (10,000 rows)

7. Conclusion

This project demonstrates a production-style AI analytics system by combining:

- LLM-based reasoning for planning and explanation
- Deterministic Pandas execution for reliability
- Modular agent design for maintainability

By enforcing deterministic execution and declarative experiment definitions, the system provides strong reproducibility guarantees and aligns with modern AI-powered business intelligence practices.