



DATA WAREHOUSE MINING (CT-463)

Department of Computer Science and Information Technology

PROJECT: DATA WAREHOUSE DESIGN IMPLEMENTATION

GROUP MEMBERS:

1. Muhammad Abdul Rafay Shaikh (CT-21035).
2. Rohan Ahmed Siddiqui (CT-21050).
3. Muhammad Umair Shaikh (CT-21048).
4. Syed Sumam Zaidi (CT-21038).

Section: A.

Contents

DATA FROM CSV	3
Introduction	3
Data Extraction:	3
Data Transformation:.....	3
Data Load:	4
Machine Learning Model	4
DATA FROM WEB SCRAPING.....	8
Introduction	8
Data Extraction	8
Data Transformation.....	8
Data Load	9
DATA FROM API.....	12
Introduction.....	12
Data Extraction	12
Data Transformation	13
Data Load.....	14

DATA FROM CSV

Introduction

The objective of this dataset is to analyze data collected from a survey (likely on employment, educational level, coding experience, and salary) and develop a machine learning model to predict annual salary. The dataset, structured with variables like Country, Education Level (EdLevel), Years of Professional Coding Experience (YearsCodePro), Employment, and Annual Compensation (ConvertedCompYearly), is explored and refined to extract insights and prepare it for predictive modeling.

Data Extraction:

The initial data extraction step involves reading the survey data into a pandas DataFrame. The data is loaded from a CSV file containing various columns, with only the relevant ones—Country, EdLevel, YearsCodePro, Employment, and ConvertedCompYearly—being retained for analysis.

```
df.head()
```

ResponseId	MainBranch	Age	Employment	RemoteWork	Check	CodingActivities	EdLevel	LearnCode	LearnCodeOnline	...	JobSatPoints_6	JobSatPoints_7	JobSatPoints_8	JobSatPoints_9	JobSatPoints_10	JobSatPoints_1
0	1	I am a developer by profession Under 18 years old	Employed, full-time	Remote	Apples	Hobby	Primary/elementary school	Books / Physical media	NaN	...	NaN	NaN	NaN	NaN	NaN	Na
1	2	I am a developer by profession 35-44 years old	Employed, full-time	Remote	Apples	Hobby;Contribute to open-source projects;Other...	Bachelor's degree (B.A., B.S., B.Eng., etc.)	Books / Physical media;Colleague;On the job tr...	Technical documentation;Blogs;Books;Written Tu...	...	0.0	0.0	0.0	0.0	0.0	0.
2	3	I am a developer by profession 45-54 years old	Employed, full-time	Remote	Apples	Hobby;Contribute to open-source projects;Other...	Master's degree (M.A., M.S., M.Eng., MBA, etc.)	Books / Physical media;Colleague;On the job tr...	Technical documentation;Blogs;Books;Written Tu...	...	NaN	NaN	NaN	NaN	NaN	Na
3	4	I am learning to code 18-24 years old	Student, full-time	NaN	Apples	NaN	Some college/university study without earning ...	Other online resources (e.g., videos, blogs, f...	Stack Overflow;How-to videos;Interactive tutorial	...	NaN	NaN	NaN	NaN	NaN	Na
4	5	I am a developer by profession 18-24 years old	Student, full-time	NaN	Apples	NaN	Secondary school (e.g. American high school, G...	Other online resources (e.g., videos, blogs, f...	Technical documentation;Blogs;Written Tutorial...	...	NaN	NaN	NaN	NaN	NaN	Na

5 rows x 114 columns

Data Transformation:

The dataset undergoes a series of transformations to ensure it is suitable for modeling. This process includes:

- **Data Cleaning:** Handling missing values, if any, in critical columns, and filtering out or imputing invalid or outlier entries.
- **Feature Engineering:** Transforming `YearsCodePro` (years of professional experience) and `EdLevel` (educational level) into numerical forms if necessary. These variables may also be normalized or categorized to improve the model's accuracy.
- **Encoding Categorical Variables:** Variables like `Country` and `Employment` are converted into numerical values through encoding methods (e.g., one-hot encoding or label encoding).
- **Data Splitting:** The data is split into training and test sets to allow for model training and evaluation.

	Country	EdLevel	YearsCodePro	Employment	Salary
72	Pakistan	Secondary school (e.g. American high school, G...	1	Employed, full-time;Student, full-time;Indepen...	7322.0
374	Austria	Professional degree (JD, MD, Ph.D, Ed.D, etc.)	6	Employed, full-time	30074.0
379	Turkey	Master's degree (M.A., M.S., M.Eng., MBA, etc.)	6	Employed, full-time	91295.0
385	France	Master's degree (M.A., M.S., M.Eng., MBA, etc.)	17	Independent contractor, freelancer, or self-em...	53703.0
389	United States of America	Some college/university study without earning ...	7	Employed, full-time;Student, part-time	110000.0

Once transformed, the prepared data is loaded into the machine learning pipeline, where it undergoes further processing, including scaling and any necessary adjustments to optimize the data for model training.

```
print("Data successfully stored in MongoDB!")
```

DATA-FROM-CSV

Data Services

Charts

+ Create Database

Q Search Namespaces

ATARIASE

DataWarehouse-Project

Developers-Data-For...

ERVICES

Search

m Processing

mx

tion

Federation

mmatic Access

ECURITY

estort

up

base Access

ork/Access

now1

DataWarehouse-Project.Developers-Data-For-Salary-Prediction

STORAGE SIZE: 468 LOGICAL DATA SIZE: 14048 TOTAL DOCUMENTS: 2016 INDEXES TOTAL SIZE: 468

Find Indexes Schema Anti-Patterns Aggregation Search Indexes

Generate queries from natural language in Composio

INSERT DOCUMENT

Filter

Type a query in: { field: "value" }

Reset Apply Options

QUERY RESULTS: 1-20 OF MANY

_id: ObjectId('672b4a854ec325b9673045')

Country: "United Kingdom of Great Britain and Northern Ireland"

EducationLevel: "Post grad"

YearsCodePro: 18

Salary: 161048

_id: ObjectId('672b4a854ec325b9673046')

Country: "United Kingdom of Great Britain and Northern Ireland"

EducationLevel: "Master's degree"

YearsCodePro: 25

Salary: 121018

_id: ObjectId('672b4a854ec325b967294f')

Country: "United States of America"

< PREVIOUS

1-20 of many results

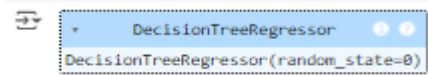
NEXT >

A machine learning model (possibly a regression model, given the focus on salary prediction) is implemented to predict the ConvertedCompYearly variable. The process includes:

1. **Model Selection:** Selection of an appropriate algorithm, potentially involving a comparison of several algorithms to determine the best fit.
2. **Training:** Training the model on the prepared dataset.

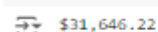
3. **Evaluation:** Evaluating model performance using standard metrics such as Mean Absolute Error (MAE) or Mean Squared Error (MSE) to gauge accuracy in salary prediction.
4. **Hyper parameter Tuning:** Fine-tuning model parameters to improve performance.

```
[ ] from sklearn.tree import DecisionTreeRegressor
    dec_tree_reg = DecisionTreeRegressor(random_state=0)
    dec_tree_reg.fit(X, y.values)
```

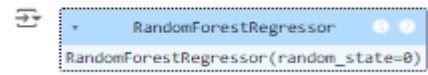


```
[ ] y_pred = dec_tree_reg.predict(X)
```

```
[ ] error = np.sqrt(mean_squared_error(y, y_pred))
    print("${:,.02f}".format(error))
```

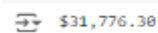


```
[ ] from sklearn.ensemble import RandomForestRegressor
    random_forest_reg = RandomForestRegressor(random_state=0)
    random_forest_reg.fit(X, y.values)
```



```
[ ] y_pred = random_forest_reg.predict(X)
```

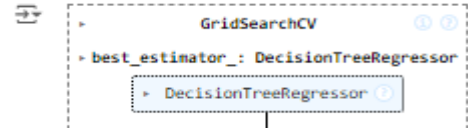
```
[ ] error = np.sqrt(mean_squared_error(y, y_pred))
    print("${:,.02f}".format(error))
```



```
[ ] from sklearn.model_selection import GridSearchCV

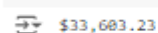
    max_depth = [None, 2,4,6,8,10,12]
    parameters = {"max_depth": max_depth}

    regressor = DecisionTreeRegressor(random_state=0)
    gs = GridSearchCV(regressor, parameters, scoring='neg_mean_squared_error')
    gs.fit(X, y.values)
```



```
[ ] regressor = gs.best_estimator_

    regressor.fit(X, y.values)
    y_pred = regressor.predict(X)
    error = np.sqrt(mean_squared_error(y, y_pred))
    print("${:,.02f}".format(error))
```



Models Tested

1. Decision Tree Regressor:

- This model partitions the data recursively based on feature values, creating a tree structure. It's known for its interpretability and ability to capture non-linear relationships.
- **RMSE:** \$31,646.22

2. Random Forest Regressor:

- An ensemble model that creates multiple decision trees and averages their predictions, which generally improves accuracy and reduces overfitting.
- **RMSE:** \$31,776.30

3. Decision Tree Regressor (Hyperparameter Tuning):

- Using **GridSearchCV**, we tuned the `max_depth` parameter of the Decision Tree Regressor to find the best-performing tree depth.
- **RMSE:** \$33,683.23

Best Model Selection

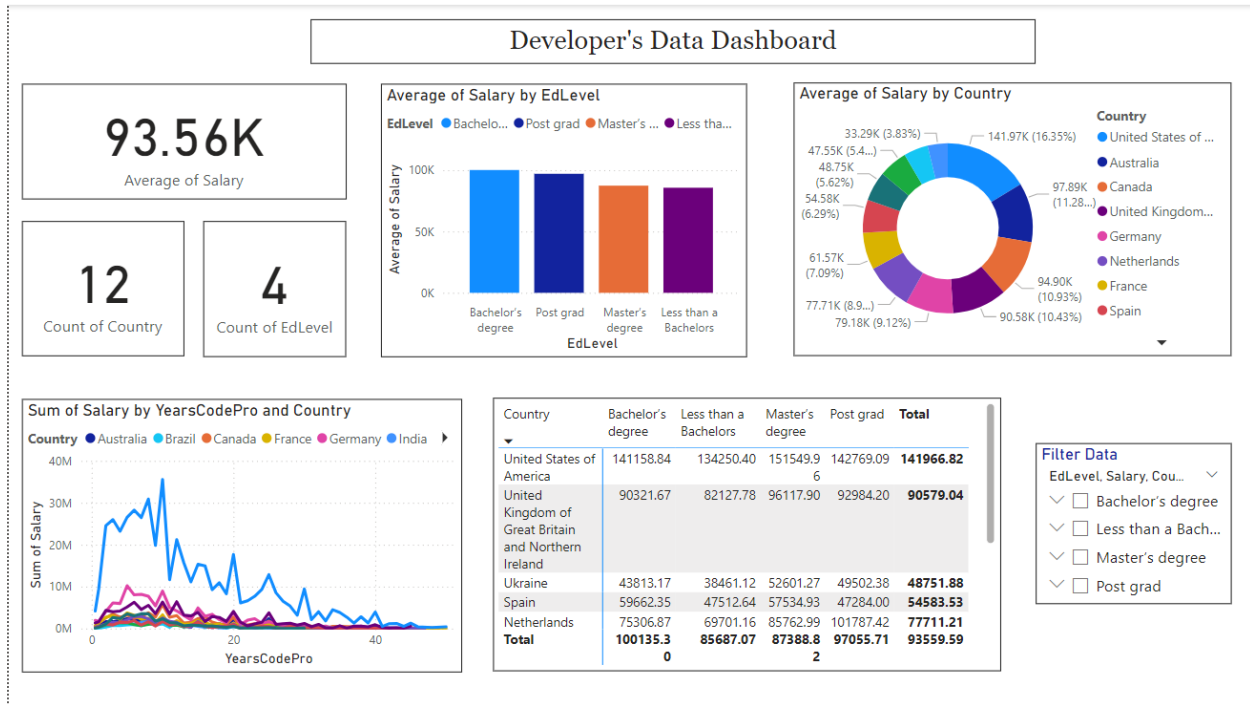
Based on the RMSE values, the **Decision Tree Regressor** without hyperparameter tuning had the lowest error, with an RMSE of \$31,646.22, making it the most accurate model for our dataset. Although the Random Forest Regressor usually performs well, in this case, the simpler Decision Tree Regressor performed slightly better.

Predictions

We used the **Decision Tree Regressor** model to make predictions on the test dataset. Below are the predictions

Example-1	Example-2
<pre>[] # country, edlevel, yearscode X = np.array([["United States", 'Master's degree', 15]]) X</pre> <pre>array([["United States", 'Master's degree', '15']], dtype='<U21') [] # Fit the encoders on all unique values in the columns le_country.fit(X[:, 0]) le_education.fit(X[:, 1]) # Transform the columns X[:, 0] = le_country.transform(X[:, 0]) X[:, 1] = le_education.transform(X[:, 1]) # Convert to float if needed X = X.astype(float) X</pre> <pre>array([[0., 0., 15.]]) [] y_pred = regressor.predict(X) y_pred</pre> <pre>/usr/local/lib/python3.10/dist-packages/sklearn/base.py:493: UserWarning: array([[101129.91304348]])</pre>	<pre>[] # country, edlevel, yearscode X = np.array([["United States", 'Master's degree', 5]]) X</pre> <pre>array([["United States", 'Master's degree', '5']], dtype='<U21') [] # Fit the encoders on all unique values in the columns le_country.fit(X[:, 0]) le_education.fit(X[:, 1]) # Transform the columns X[:, 0] = le_country.transform(X[:, 0]) X[:, 1] = le_education.transform(X[:, 1]) # Convert to float if needed X = X.astype(float) X</pre> <pre>array([[0., 0., 5.]]) [] y_pred = regressor.predict(X) y_pred</pre> <pre>/usr/local/lib/python3.10/dist-packages/sklearn/base.py:493: UserWarning: array([[84275.5]])</pre>

With the help of PowerBI, Dashboard has been created to analyze insights.



DATA FROM WEB SCRAPING

Introduction

The purpose of this project is to predict car prices based on data scraped from the PakWheels website. PakWheels is a popular platform for buying and selling vehicles, and it provides valuable insights into car market trends in Pakistan. This project aims to gather data from PakWheels, process it, store it in a database, and then use machine learning models to predict car prices based on features such as model, year, and other car specifications.

Data Extraction

To collect data from PakWheels, we used web scraping techniques to gather information from car listings. This process involved:

- **HTML Tags:** We navigated through the HTML structure of the PakWheels website to locate relevant data elements. Each car listing had specific tags for information like model, year, price, and contact details. By examining these tags, we identified where each piece of data was stored within the HTML document.
- **strip() Function:** The `strip()` function was used to clean up the text data extracted from these tags, removing unnecessary whitespace and formatting issues. This helped in making the data consistent and ready for further processing.

This extraction process resulted in a dataset with essential car details, either showing a price directly or indicating that the user should call for price information.

Data Transformation

After extracting the raw data from PakWheels, several data transformation steps were performed to clean and prepare the dataset for analysis and modeling. This process included **handling missing values**, **splitting the data into subsets**, and **applying label encoding** to categorical variables.

Handling Missing Values

The dataset contained some missing values, especially in fields like price (when listings prompted users to call for the price) and certain car specifications. To manage this:

1. **Price Column:** Listings without explicit price values were categorized separately into the `data_call` subset (listings with "call" instead of a price). Listings with price values were placed into the `data_update` subset, which was used for predictive modeling.
2. **Other Missing Values:** For missing values in other columns, such as car specifications, we either imputed these values with a suitable placeholder (such as "Unknown") or dropped rows with extensive missing data, depending on the importance of the feature for our analysis.

This process ensured that the data was clean and ready for machine learning, with minimal missing information that could interfere with model training.

Label Encoding

To enable machine learning algorithms to process categorical data, we applied **label encoding** to relevant categorical features. Label encoding is a technique that converts categorical variables (text) into numerical values, allowing models to interpret them effectively.

1. **Car Model and Brand:** Car models and brands were label-encoded so that each unique model and brand was represented by a unique numerical identifier.
2. **Fuel Type, Transmission, and Condition:** Other categorical fields, such as fuel type (e.g., Petrol, Diesel), transmission (e.g., Manual, Automatic), and car condition (e.g., New, Used), were also label-encoded.

Label encoding helped to transform the data into a numerical format, making it compatible with the machine learning algorithms while preserving the distinct categories.

```
✓ [20] data_update.head()
```

	Manufacturer	Variant	Details	Location	Model	Distance Travelled	Fuel Type	Engine Capacity	Transmission	Featured	Price in PKR
0	29	103	819	175	2020	72000	5	1800.0	1	1	5800000.0
1	64	295	1104	140	2010	139000	5	1000.0	1	1	2100000.0
2	29	103	819	109	2019	52000	5	1800.0	1	1	6050000.0
3	29	102	343	140	2019	64000	5	1500.0	1	1	4125000.0
4	64	327	718	66	2022	18500	3	1200.0	1	1	9000000.0

Data Splitting

Following data cleaning and encoding, we split the dataset into two main subsets:

1. **data_call:** This subset contained listings where users were prompted to "call" for price information. Since it did not contain actual price values, this subset was not used for model training but could be retained for future analysis.
2. **data_update:** This subset contained listings with explicit price values. It served as the primary dataset for machine learning, as it provided the target variable (`price`) necessary for model training.

Data Load

The transformed data was stored in MongoDB for easy access and scalability. We created two separate **collections** in MongoDB to store each subset:

```

✓ [34] # Insert data_call into the 'Used_Cars_call' collection
18s db['Used_Cars_call'].insert_many(data_call.to_dict("records"))

# Insert data_update into the 'Used_Cars_pkr' collection
db['Used_Cars_pkr'].insert_many(data_update.to_dict("records"))

11s ObjectId('672de9546bf58e5986a43552'), ObjectId('672de9546bf58e5986a43553'), ObjectId('672de9546bf58e5986a43554'), ObjectId('672de9546bf58e5986a43555'), ObjectId('672de9546bf58e5986a43556'),
ObjectId('672de9546bf58e5986a43557'), ObjectId('672de9546bf58e5986a43558'), ObjectId('672de9546bf58e5986a43559'), ObjectId('672de9546bf58e5986a4355a'), ObjectId('672de9546bf58e5986a4355b'),
ObjectId('672de9546bf58e5986a4355c'), ObjectId('672de9546bf58e5986a4355d'), ObjectId('672de9546bf58e5986a4355e'), ObjectId('672de9546bf58e5986a4355f'), ObjectId('672de9546bf58e5986a43560'),
ObjectId('672de9546bf58e5986a43561'), ObjectId('672de9546bf58e5986a43562'), ObjectId('672de9546bf58e5986a43563'), ObjectId('672de9546bf58e5986a43564'), ObjectId('672de9546bf58e5986a43565'),
ObjectId('672de9546bf58e5986a43566'), ObjectId('672de9546bf58e5986a43567'), ObjectId('672de9546bf58e5986a43568'), ObjectId('672de9546bf58e5986a43569'), ObjectId('672de9546bf58e5986a4356a'),
ObjectId('672de9546bf58e5986a4356b'), ObjectId('672de9546bf58e5986a4356c'), ObjectId('672de9546bf58e5986a4356d'), ObjectId('672de9546bf58e5986a4356e'), ObjectId('672de9546bf58e5986a4356f'),
ObjectId('672de9546bf58e5986a43570'), ObjectId('672de9546bf58e5986a43571'), ObjectId('672de9546bf58e5986a43572'), ObjectId('672de9546bf58e5986a43573'), ObjectId('672de9546bf58e5986a43574'),
ObjectId('672de9546bf58e5986a43575'), ObjectId('672de9546bf58e5986a43576'), ObjectId('672de9546bf58e5986a43577'), ObjectId('672de9546bf58e5986a43578'), ObjectId('672de9546bf58e5986a43579'),
ObjectId('672de9546bf58e5986a4357a'), ObjectId('672de9546bf58e5986a4357b'), ObjectId('672de9546bf58e5986a4357c'), ObjectId('672de9546bf58e5986a4357d'), ObjectId('672de9546bf58e5986a4357e'),
ObjectId('672de9546bf58e5986a4357f'), ObjectId('672de9546bf58e5986a43580'), ObjectId('672de9546bf58e5986a43581'), ObjectId('672de9546bf58e5986a43582'), ObjectId('672de9546bf58e5986a43583'),
ObjectId('672de9546bf58e5986a43584'), ObjectId('672de9546bf58e5986a43585'), ObjectId('672de9546bf58e5986a43586'), ObjectId('672de9546bf58e5986a43587'), ObjectId('672de9546bf58e5986a43588'),
ObjectId('672de9546bf58e5986a43589'), ObjectId('672de9546bf58e5986a4358a'), ObjectId('672de9546bf58e5986a4358b'), ObjectId('672de9546bf58e5986a4358c'), ObjectId('672de9546bf58e5986a4358d'),
ObjectId('672de9546bf58e5986a4358e'), ObjectId('672de9546bf58e5986a4358f'), ObjectId('672de9546bf58e5986a43590'), ObjectId('672de9546bf58e5986a43591'), ObjectId('672de9546bf58e5986a43592'),
ObjectId('672de9546bf58e5986a43593'), ObjectId('672de9546bf58e5986a43594'), ObjectId('672de9546bf58e5986a43595'), ObjectId('672de9546bf58e5986a43596'), ObjectId('672de9546bf58e5986a43597'),
ObjectId('672de9546bf58e5986a43598'), ObjectId('672de9546bf58e5986a43599'), ObjectId('672de9546bf58e5986a4359a'), ObjectId('672de9546bf58e5986a4359b'), ObjectId('672de9546bf58e5986a4359c'),
ObjectId('672de9546bf58e5986a4359d'), ObjectId('672de9546bf58e5986a4359e'), ObjectId('672de9546bf58e5986a4359f'), ObjectId('672de9546bf58e5986a435a0'), ObjectId('672de9546bf58e5986a435a1'),
ObjectId('672de9546bf58e5986a435a2'), ObjectId('672de9546bf58e5986a435a3'), ObjectId('672de9546bf58e5986a435a4'), ObjectId('672de9546bf58e5986a435a5'), ObjectId('672de9546bf58e5986a435a6'),
ObjectId('672de9546bf58e5986a435a7'), ObjectId('672de9546bf58e5986a435a8'), ObjectId('672de9546bf58e5986a435a9'), ObjectId('672de9546bf58e5986a435aa'), ObjectId('672de9546bf58e5986a435ab'),
ObjectId('672de9546bf58e5986a435ac'), ObjectId('672de9546bf58e5986a435ad'), ObjectId('672de9546bf58e5986a435ae'), ObjectId('672de9546bf58e5986a435af'), ObjectId('672de9546bf58e5986a435b0'),
ObjectId('672de9546bf58e5986a435b1'), ObjectId('672de9546bf58e5986a435b2'), ObjectId('672de9546bf58e5986a435b3'), ObjectId('672de9546bf58e5986a435b4'), ObjectId('672de9546bf58e5986a435b5'),
ObjectId('672de9546bf58e5986a435b6'), ObjectId('672de9546bf58e5986a435b7'), ObjectId('672de9546bf58e5986a435b8'), ObjectId('672de9546bf58e5986a435b9'), ObjectId('672de9546bf58e5986a435ba'),
ObjectId('672de9546bf58e5986a435bb'), ObjectId('672de9546bf58e5986a435bc'), ObjectId('672de9546bf58e5986a435bd'), ObjectId('672de9546bf58e5986a435be'), ObjectId('672de9546bf58e5986a435bf'),
ObjectId('672de9546bf58e5986a435c0'), ObjectId('672de9546bf58e5986a435c1'), ObjectId('672de9546bf58e5986a435c2'), ObjectId('672de9546bf58e5986a435c3'), ObjectId('672de9546bf58e5986a435c4'),
ObjectId('672de9546bf58e5986a435c5'), ObjectId('672de9546bf58e5986a435c6'), ObjectId('672de9546bf58e5986a435c7'), ObjectId('672de9546bf58e5986a435c8'), ObjectId('672de9546bf58e5986a435c9'),
ObjectId('672de9546bf58e5986a435ca'), ObjectId('672de9546bf58e5986a435cb'), ObjectId('672de9546bf58e5986a435cc'), ObjectId('672de9546bf58e5986a435cd'), ObjectId('672de9546bf58e5986a435ce'),
ObjectId('672de9546bf58e5986a435cf'), ObjectId('672de9546bf58e5986a435d0'), ObjectId('672de9546bf58e5986a435d1'), ObjectId('672de9546bf58e5986a435d2'), ObjectId('672de9546bf58e5986a435d3'),
ObjectId('672de9546bf58e5986a435d4'), ObjectId('672de9546bf58e5986a435d5'), ObjectId('672de9546bf58e5986a435d6'), ObjectId('672de9546bf58e5986a435d7'), ObjectId('672de9546bf58e5986a435d8'),
ObjectId('672de9546bf58e5986a435d9'), ObjectId('672de9546bf58e5986a435da'), ObjectId('672de9546bf58e5986a435db'), ObjectId('672de9546bf58e5986a435dc'), ObjectId('672de9546bf58e5986a435dd'),
ObjectId('672de9546bf58e5986a435de'), ObjectId('672de9546bf58e5986a435df'), ObjectId('672de9546bf58e5986a435e0'), ObjectId('672de9546bf58e5986a435e1'), ObjectId('672de9546bf58e5986a435e2'),
ObjectId('672de9546bf58e5986a435e3'), ObjectId('672de9546bf58e5986a435e4'), ObjectId('672de9546bf58e5986a435e5'), ObjectId('672de9546bf58e5986a435e6'), ObjectId('672de9546bf58e5986a435e7'),
ObjectId('672de9546bf58e5986a435e8'), ObjectId('672de9546bf58e5986a435e9'), ObjectId('672de9546bf58e5986a435ea'), ObjectId('672de9546bf58e5986a435eb'), ObjectId('672de9546bf58e5986a435ec'),
ObjectId('672de9546bf58e5986a435ed'), ObjectId('672de9546bf58e5986a435ee'), ObjectId('672de9546bf58e5986a435ef'), ObjectId('672de9546bf58e5986a435f0'), ObjectId('672de9546bf58e5986a435f1'),
ObjectId('672de9546bf58e5986a435f2'), ObjectId('672de9546bf58e5986a435f3'), ObjectId('672de9546bf58e5986a435f4'), ObjectId('672de9546bf58e5986a435f5'), ObjectId('672de9546bf58e5986a435f6'),

```

- **data_call collection:** This collection contains the listings where users are prompted to call for price details.
- **data_update collection:** This collection stores listings with an explicit price, making it the primary dataset for our machine learning model.

Create Database

Search Homepages

DataWarehouse-Project

Covid-19-Data

Developers-Datas-For-S...

Used_Cars_call

Used_Cars_pkr

DataWarehouse-Project-Used_Cars_pkr

STORAGE SIZE: 24MB

LOGICAL DATA SIZE: 182MB

TOTAL DOCUMENTS: 4849

INDEXED TOTAL SIZE: 24MB

Find Indexes Schema And Patterns () Aggregation Search Indexes

Generate queries from natural language in Compass!

File #

Type a query: (Field: 'value')

Reset Apply Options

QUERY RESULTS: 120 OF MANY

_id: ObjectId('672de9546bf58e5986a435c0')

Manufacturer: 29

Variant: 281

Seatfix: 639

Location: 175

Model: 1038

Distance Traveled: 12084

Fuel Type: 1

Engine Capacity: 1386

Transmission: 1

Featured: 1

Price In PKR: 300000

_id: ObjectId('672de9546bf58e5986a435d4')

Manufacturer: 94

Variant: 281

PREVIOUS

120 of many results

NEXT

Create Database

Search Homepages

DataWarehouse-Project

Covid-19-Data

Developers-Datas-For-S...

Used_Cars_call

Used_Cars_pkr

DataWarehouse-Project-Used_Cars_call

STORAGE SIZE: 24MB

LOGICAL DATA SIZE: 174MB

TOTAL DOCUMENTS: 2044

INDEXED TOTAL SIZE: 24MB

Find Indexes Schema And Patterns () Aggregation Search Indexes

Generate queries from natural language in Compass!

File #

Type a query: (Field: 'value')

Reset Apply Options

QUERY RESULTS: 120 OF MANY

_id: ObjectId('672de9546bf58e5986a435c7')

Manufacturer: 1994

Variant: 848

Seatfix: 7031

Location: 1514mbp

Model: 1832

Distance Traveled: 10088

Fuel Type: Petrol

Engine Capacity: 1280

Transmission: Automatic

Featured: 104

Price In PKR: 11111

_id: ObjectId('672de9546bf58e5986a435e4')

Manufacturer: 1994

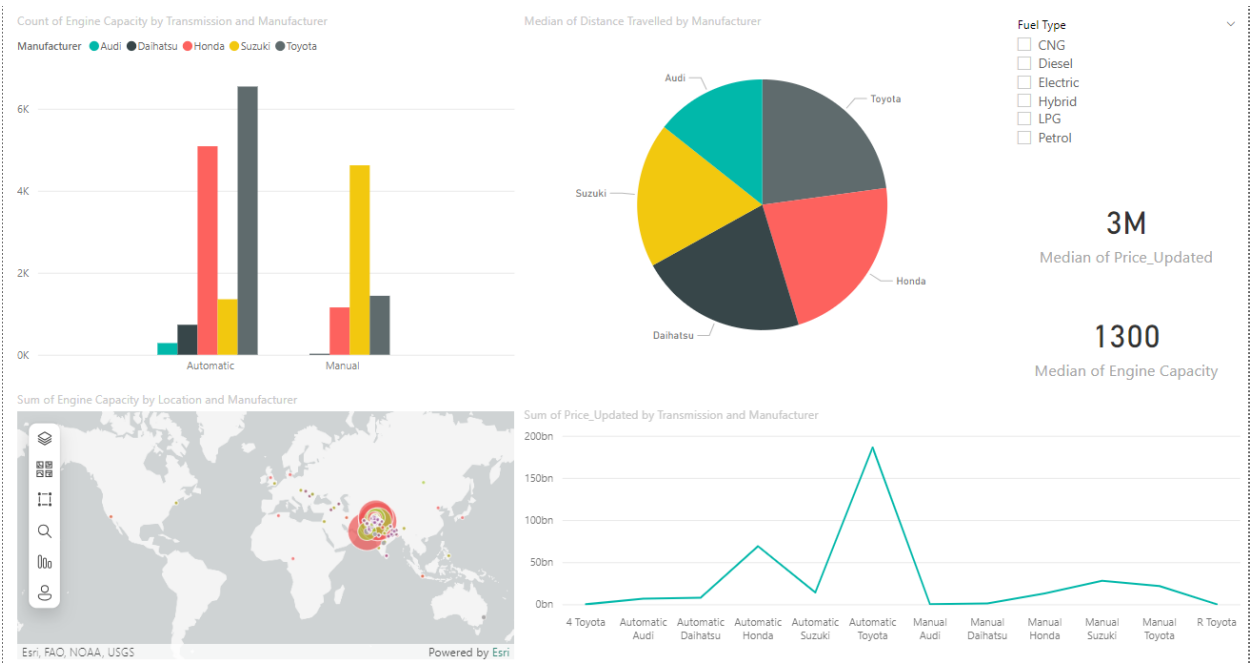
Variant: 848

PREVIOUS

120 of many results

NEXT

Dashboard



DATA FROM API

Introduction

In today's data-driven landscape, centralized data repositories are crucial for deriving meaningful insights. This project involves creating a data warehouse using the ETL (Extract, Transform, Load) process to consolidate COVID-19 statistics from an external API. By extracting, transforming, and loading this data into MongoDB, we facilitate easy access and analysis of pandemic-related information.

Data Extraction

Data extraction is the first step in the ETL pipeline, where raw data is sourced from the COVID-19 API. This API provides comprehensive statistics on COVID-19, such as infection rates, recoveries, and testing figures by country.

Libraries Used: The requests library is used to make API calls, and json is utilized for parsing the API's JSON response, allowing smooth data retrieval.

```
covid = requests.get(url, headers=headers).json()

print(json.dumps(covid, indent=4))

{
  "get": "statistics",
  "parameters": [],
  "errors": [],
  "results": 238,
  "response": [
    {
      "continent": "Africa",
      "country": "Djibouti",
      "population": 1016097,
      "cases": {
        "new": null,
        "active": 74,
        "critical": null,
        "recovered": 15427,
        "1M_pop": "15441",
        "total": 15690
      },
      "deaths": {
        "new": null,
        "1M_pop": "186",
        "total": 189
      },
      "tests": {
        "1M_pop": "301094",
        "total": 305941
      },
      "day": "2024-11-06",
      "time": "2024-11-06T06:45:09+00:00"
    },
  ],
}
```

API URL: <https://covid-193.p.rapidapi.com/statistics>

Data Transformation

The transformation stage ensures that data is standardized, cleaned, and optimized for storage and querying. Key transformations include:

Data Cleaning: Using Pandas, missing values are handled, data types are corrected, and country names are standardized.

Schema Structuring: Data is organized into a structured format that aligns with MongoDB's document-based structure.

Libraries Used: Pandas for data manipulation and NumPy for numerical operations to ensure efficient data processing.

```
[ ] covid_data
```



	Country	Population	Recovered	TotalTests
0	China	1.448471e+09	379053.0	1.600000e+08
1	Falkland-Islands	3.539000e+03	1930.0	8.632000e+03
2	Montserrat	4.965000e+03	1376.0	1.776200e+04
3	Nauru	1.090300e+04	5347.0	2.050900e+04
4	Wallis-and-Futuna	1.098200e+04	438.0	2.050800e+04
...
165	USA	3.348053e+08	109814428.0	1.186852e+09
166	France	6.558452e+07	39970918.0	2.714902e+08
167	Germany	8.388360e+07	38240600.0	1.223324e+08
168	Brazil	2.153536e+08	36249161.0	6.377617e+07
169	S-Korea	5.132990e+07	34535939.0	1.580406e+07

170 rows x 4 columns

```
[ ] pakistan_data = covid_data[covid_data['Country'] == 'Pakistan']
```

```
[ ] pakistan_data
```



	Country	Population	Recovered	TotalTests
120	Pakistan	229488994.0	1538689.0	30589153.0

Data Load

In the final ETL phase, the transformed data is loaded into a MongoDB database. MongoDB's document-oriented model is well-suited for storing JSON-like data, making it ideal for handling the COVID-19 statistics.

Database: Data is stored in a MongoDB collection, providing a flexible and scalable storage solution.

Batch Loading: Data is inserted into MongoDB in batches, ensuring a smooth loading process and efficient access.

Libraries Used: pymongo for connecting to MongoDB and inserting data into collections

```
[ ] data_dict = covid_data.to_dict("records") # Convert DataFrame to a list of dictionaries
collection.insert_many(data_dict)
```

The screenshot displays the MongoDB Atlas web interface. On the left, a sidebar contains navigation links: Overview, Clusters, SERVICES, Atlas Search, Stream Processing, Triggers, Migration, Data Federation, Programmatic Access, SECURITY, Quickstart, Backup, Database Access, Network Access, Advanced, and Oato. The main content area is titled 'DataWarehouse-Project.Covid-API-Data'. It shows the collection's metadata: STORAGE SIZE: 28KB, LOGICAL DATA SIZE: 156KB, TOTAL DOCUMENTS: 170, and INDEXES TOTAL SIZE: 20KB. Below this, there are tabs for Find, Indexes, Schema Anti-Patterns, Aggregation, and Search Indexes. A 'Generate queries from natural language in Compass' link is present. A 'Filter' section allows users to type a query (e.g., 'field: 'value''). Below the filter, 'QUERY RESULTS 1-20 OF MANY' are displayed. The first three results are visible: 1. _id: Object('672b53d1579f6f3f538f'), Country: 'Falkland-Islands', Population: 3539, Recovered: 1930, TotalTests: 8632. 2. _id: Object('672b53d1579f6f3f538e'), Country: 'Montserrat', Population: 4965, Recovered: 1336, TotalTests: 13762. 3. _id: Object('672b53d1579f6f3f538f'), Country: 'China'. At the bottom, there are pagination controls showing '1-20 of many results' and 'NEXT' button.