

---

---

# Global Carbon Emission Comparisons

Group 15 - Ruchita, Avneesh, Tarini, Kobe, Mikhail

---

---

# Topic Description

- Alarming rates of carbon emissions
  - Study major factors affecting climate conditions
  - Analyze correlations between energy usage and GHG emissions through transportation
  - Determine the recent impact of alternative energy sources
  - Predict future impact of higher alternative energy adoption
-

---

# Group Roles

- X - Avneesh
  - Triangle - Tarini
  - Circle - Mikhail
  - Square - Ruchita & Kobe
-

---

# Data Description

- Annual worldwide vehicle mix by different energy types (Diesel, Electric, Hybrid, etc)
  - Overall GHG emissions by sector
  - Market demand of all vehicles types
  - Energy production
  - Gathered data to determine our main goal of proving that the introduction of EV will decrease annual gas emissions on a global scale
-

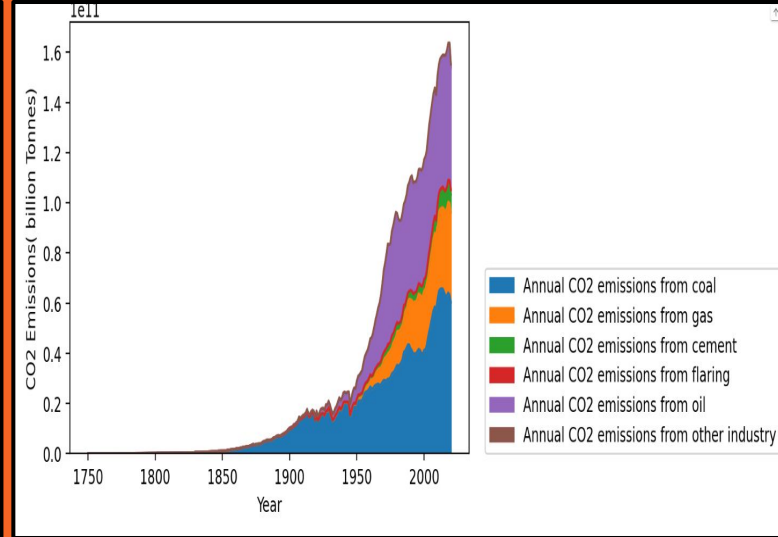
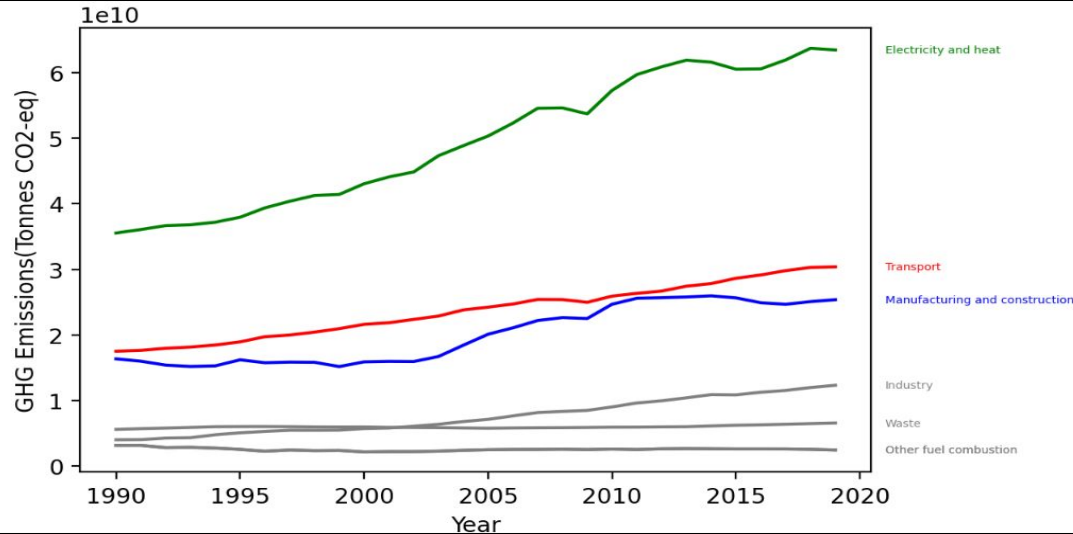
---

# Machine Learning Preliminary Process

*Conducted Data Exploration to cover the following:*

- All datasets were merged based on country and year
  - Extracted EV demand dataset from original sheet
  - Remove non-value variables and refined the column names to match across all merged datasets
  - Loaded the GHG emissions and EV demand datasets as pandas dataframe
  - Defined the merge logic to join both the datasets
-

# Data Analysis Phase



- The major focus was to analyze how much CO<sub>2</sub> emissions were saved when introducing EV variable into the transport sector.
- Out of all factors, **electricity & heat**, **transport**, and **manufacturing & construction**.

---

# Database

- Downloaded datasets related to co2 emissions and GHG emissions
  - Source - <https://ourworldindata.org/>
- Tried to find the unique key to merge the data in this case it is the country name and year
- Collected the datasets for Transport Vehicle demand country wise and year wise, merged it with the GHG emissions database

# Database Sample

The screenshot displays the PostgreSQL interface with a query editor and a data output table. The query editor shows a simple SELECT statement: `select * from merged`. The data output table lists 11 rows of data from the 'merged' table, including columns for a primary key 'num', 'Entity', 'Year', and various electricity-related metrics.

	num [PK] integer	Entity character varying	Year integer	Fossil fuels (% electricity) double precision	Renewables (% electricity) double precision	Per capita electricity (kWh) double precision	Fossil fuels double prec
1	0	Africa	2000	78.74517059	17.65452576	542.7769165	
2	1	Africa	2001	79.02587128	18.2166481	552.5474243	
3	2	Africa	2002	79.12548065	18.00189972	572.7458496	
4	3	Africa	2003	80.06246948	16.8457756	583.6015625	
5	4	Africa	2004	80.04415894	16.82542038	605.0234375	
6	5	Africa	2005	80.85244751	16.53281784	611.7351685	
7	6	Africa	2006	81.17703247	16.49860191	624.1119995	
8	7	Africa	2007	81.25521088	16.27716827	634.1665649	
9	8	Africa	2008	81.09965515	16.29686546	631.3635864	
10	9	Africa	2009	80.7677536	16.63348198	621.9937744	
11	10	Africa	2010	80.49606900	16.0411960	647.7619400	



# Machine Learning Model Summary

## OLS Regression Results

Dep. Variable:	ev_savings	R-squared (uncentered):	0.507			
Model:	OLS	Adj. R-squared (uncentered):	0.493			
Method:	Least Squares	F-statistic:	36.02			
Date:	Thu, 01 Sep 2022	Prob (F-statistic):	2.36e-75			
Time:	03:50:38	Log-Likelihood:	-8982.0			
No. Observations:	576	AIC:	1.800e+04			
Df Residuals:	560	BIC:	1.807e+04			
Df Model:	16					
Covariance Type: nonrobust						
	coef	std err	t	P> t	[0.025	0.975]
Fossil fuels (% equivalent primary energy)	7950.8993	5231.081	1.520	0.129	-2324.038	1.82e+04
Transport_x	0.1986	0.161	1.234	0.218	-0.118	0.515
Agriculture	0.0011	0.001	1.013	0.312	-0.001	0.003
Land-use change and forestry	0.0015	0.001	2.731	0.007	0.000	0.003
Waste	0.0082	0.005	1.731	0.084	-0.001	0.018
Industry	0.0584	0.005	10.943	0.000	0.048	0.069
Manufacturing and construction	-0.0240	0.002	-11.381	0.000	-0.028	-0.020
Transport_y	-0.1989	0.155	-1.280	0.201	-0.504	0.106
Electricity and heat	0.0080	0.001	7.096	0.000	0.006	0.010
Buildings	-0.0098	0.003	-3.043	0.002	-0.016	-0.003
Fugitive emissions	-0.0062	0.002	-3.892	0.000	-0.009	-0.003
Other fuel combustion	0.0040	0.016	0.246	0.805	-0.028	0.036
Aviation and shipping	-0.0209	0.006	-3.501	0.001	-0.033	-0.009
Renewables (% equivalent primary energy)	5839.8990	1.32e+04	0.444	0.657	-2e+04	3.17e+04
Fossil fuels (% electricity)	-1.103e+04	5700.997	-1.935	0.054	-2.22e+04	168.776
Renewables (% electricity)	-4951.9647	9257.701	-0.535	0.593	-2.31e+04	1.32e+04

Renewables (% equivalent primary energy)	5839.8990	1.32e+04	0.444	0.657	-2e+04	3.17e+04
Fossil fuels (% electricity)	-1.103e+04	5700.997	-1.935	0.054	-2.22e+04	168.776
Renewables (% electricity)	-4951.9647	9257.701	-0.535	0.593	-2.31e+04	1.32e+04
Omnibus:	416.999	Durbin-Watson:	1.018			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	21995.226			
Skew:	2.579	Prob(JB):	0.00			
Kurtosis:	32.831	Cond. No.	2.41e+08			

## Notes:

- [1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.
- [2] Standard Errors assume that the covariance matrix of the errors is correctly specified.
- [3] The condition number is large, 2.41e+08. This might indicate that there are strong multicollinearity or other numerical problems.

# Machine Learning Model - Data columns

```
Index([
    ('Petrol', 'Econ'),
    ('Petrol', 'Mid'),
    ('Petrol', 'Lux'),
    ('Adv Petrol', 'Econ'),
    ('Adv Petrol', 'Mid'),
    ('Adv Petrol', 'Lux'),
    ('Diesel', 'Econ'),
    ('Diesel', 'Mid'),
    ('Diesel', 'Lux'),
    ('Adv Diesel', 'Econ'),
    ('Adv Diesel', 'Mid'),
    ('Adv Diesel', 'Lux'),
    ('CNG', 'Econ'),
    ('CNG', 'Mid'),
    ('CNG', 'Lux'),
    ('Hybrid', 'Econ'),
    ('Hybrid', 'Mid'),
    ('Hybrid', 'Lux'),
    ('Electric', 'Econ'),
    ('Electric', 'Mid'),
    ('Electric', 'Lux'),
    ('Bikes', 'Econ'),
    ('Bikes', 'Lux'),
    ('Electric Bikes', 'Adv Econ'),
    ('Electric Bikes', 'Adv Lux'),
    'Fossil fuels (% equivalent primary energy)',
    'Transport_x',
    'Agriculture',
    'Land-use change and forestry',
    'Waste',
    'Industry',
    'Manufacturing and construction',
    'Transport_y',
    'Electricity and heat',
    'Buildings',
    'Fugitive emissions',
    'Other fuel combustion',
    'Aviation and shipping',
    'Renewables (% equivalent primary energy)',
    'Fossil fuels (% electricity)',
    'Renewables (% electricity)',
    'fossil_demand',
    'ev_demand',
```

- To calculate CO2 emissions saved by electric vehicles year by year
- Features given in the screenshot are taken into consideration

# Machine Learning Model Output

```
[121] X = df_ev_ghg_merged[[  
    'Fossil fuels (% equivalent primary energy)',  
    'Transport_x',  
    'Agriculture',  
    'Land-use change and forestry',  
    'Waste', 'Industry', 'Manufacturing and construction', 'Transport_y',  
    'Electricity and heat',  
    'Buildings', 'Fugitive emissions', 'Other fuel combustion', 'Aviation and shipping', 'Renewables (% equivalent primary energy)', 'Fossil fuels (% electricity)',  
    'Renewables (% electricity)']]
```

```
[122] X
```

```
▶ y = df_ev_ghg_merged[['ev_savings']]
```

```
[124] lm = linear_model.LinearRegression()  
      model = lm.fit(X,y)
```

```
▶ y_predictions = lm.predict(X)
```

```
[126] lm.score(X,y)
```

```
0.4954363469172296
```

---

# Machine Learning Model Output

```
[116] lm = linear_model.LinearRegression()
      model = lm.fit(X_train,y_train)

[117] y_predictions = lm.predict(X_test)

[118] from sklearn.metrics import r2_score

[119] r2_score(y_test, y_predictions)

0.727704844938476
```

- Coefficient of determination(r square) = 0.72 or 70%
  - This indicates that 72% variability is explained by our model.
-

# Machine Learning Data Split

- Split the original dataset into training and test split using scikit-learn and train\_test\_split module
- We train the model on train split then we test the model on X values of test split

```
Split the Data into Training and Testing

# Create our features
X = df_ev_ghg_merged[['Transport_x',
'Agriculture',
'Land-use change and forestry',
'Waste', 'Industry', 'Manufacturing and construction', 'Transport_y',
'Electricity and heat',
'Buildings', 'Fugitive emissions', 'Other fuel combustion', 'Aviation and shipping', 'Renewables (% equivalent primary energy)',
'Renewables (% electricity)']]

# Create our target
y = df_ev_ghg_merged[['ev_savings']]

[115] from sklearn.model_selection import train_test_split
      X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=1)

[116] lm = linear_model.LinearRegression()
      model = lm.fit(X_train, y_train)

[117] y_predictions = lm.predict(X_test)

[118] from sklearn.metrics import r2_score

[119] r2_score(y_test, y_predictions)

0.727704844938476
```

---

# Machine Learning Model Choice

## Linear Regression Model

### *Limitations:*

- Accuracy was initially low when including electric bike data
- For some of the features, standard error was really high, which required dropping columns

### *Benefits:*

- Accuracy increased by 20%
  - R-squared variance also increased by 20%
-

---

# Dashboard Blueprint

## *Tools:*

- PostgreSQL
- Scikit-Learn
- Jupyter Notebook
- Python
  - Matplotlib & pandas

## *Interactive Elements:*

- Graphical visualization of our findings and key takeaways.

## *Findings & Recommendations*

- Electricity and transport are the highest drivers of carbon/GHG emissions
- Selecting a country/region to determine patterns of increased usage of alternate energy
- Cleaner energy sources to aid EV charging capabilities