# Traffic flow prediction using the METR-LA traffic

## Problem Description :

The objective of this project is to predict the traffic flow using the METR-LA traffic dataset. The METR-LA traffic dataset is a public dataset that contains traffic information from loop detectors in the Los Angeles area from March 1st, 2012 to June 30th, 2012. The dataset contains traffic information from 207 sensors, with a time granularity of 5 minutes.

The goal of this project is to build a deep learning model that can accurately predict the traffic flow at a given sensor location and time. This model can be used by transportation authorities to better manage traffic flow and improve travel times for commuters.

**Dataset Description:**

The METR-LA traffic dataset contains traffic information from 207 sensors in the Los Angeles area. Each row in the dataset represents the traffic flow at a given sensor location and time, with a time granularity of 5 minutes. The dataset contains the following columns:

- **sensor_id**: The unique identifier for the sensor location.
- **date**: The date and time of the traffic flow measurement.
- **flow**: The traffic flow at the sensor location.

**Project Requirements:**

The project requirements are as follows:

1. Load and preprocess the METR-LA traffic dataset.
2. Implement a deep learning model to predict the traffic flow.
3. Train and evaluate the model using appropriate evaluation metrics.
4. Visualize the results of the model.

**Deliverables:**

The deliverables for this project are as follows:

1. A Jupyter Notebook containing the code for loading and preprocessing the METR-LA traffic dataset.
2. A Jupyter Notebook containing the code for implementing the deep learning model to predict the traffic flow.
3. A Jupyter Notebook containing the code for training and evaluating the model using appropriate evaluation metrics.
4. A Jupyter Notebook containing the code for visualizing the results of the model.

# Possible Framework:

1. **Data Loading and Preprocessing:** This section involves loading the METR-LA traffic dataset, removing missing or invalid data points, and splitting the data into training and testing sets.
2. **Feature Engineering:** This section involves engineering the features for the deep learning model. This can include time-based features such as hour of the day and day of the week, as well as sensor-specific features such as location and road type.
3. **Model Selection:** This section involves selecting the appropriate deep learning model for predicting the traffic flow. This can include models such as LSTMs, CNNs, and hybrid models.
4. **Hyperparameter Tuning:** This section involves tuning the hyperparameters of the deep learning model to optimize the performance.
5. **Model Training and Evaluation:** This section involves training the deep learning model on the training data and evaluating the model using appropriate evaluation metrics such as mean absolute error and mean squared error.
6. **Results and Visualization:** This section involves visualizing the results of the deep learning model using plots and graphs.

**Implementation Steps:**

The following are the detailed steps to implement the framework:

1. Load the METR-LA traffic dataset using pandas and preprocess the data by removing missing or invalid data points.
2. Split the data into training and testing sets.
3. Engineer the features for the deep learning model using time-based and sensor-specific features.
4. Select the appropriate deep learning model for predicting the traffic flow.
5. Tune the hyperparameters of the deep learning model using grid search or random search.
6. Train the deep learning model on the training data and evaluate the model using appropriate evaluation metrics such as mean absolute error and mean squared error.
7. Visualize the results of the deep learning model using plots and graphs.

# Code Explanation :

Here is the simple explanation for the code which is provided in the code.py file.

**Code Explanation:**

1. The first step is to load the METR-LA traffic dataset, remove missing or invalid data points, and split the data into training and testing sets.
2. The second step involves engineering the features for the deep learning model. This includes adding time-based features such as hour and weekday, as well as sensor-specific features such as rolling mean, rolling standard deviation, and difference from the previous time step.
3. The third step is to define the LSTM model for traffic flow prediction. This involves creating an LSTM layer with 64 units and a dense layer with one unit. The loss function is mean squared error, and the optimizer is Adam.
4. The fourth step is to train the LSTM model on the training data using hyperparameters such as batch size and number of epochs. After training the model, we evaluate its performance on the testing data by calculating mean squared error, mean absolute error, and R2 score.
5. The fifth step involves visualizing the actual and predicted traffic flow values on a graph.

**Requirements:**

To run the code, you need to have the following libraries installed in your Python environment:

- pandas
- numpy
- scikit-learn
- TensorFlow
- matplotlib

In addition, you need to download the METR-LA traffic dataset and save it as a CSV file named 'metr-la.h5' in the local directory.

**Running the code:**

To run the code, simply copy and paste it into a Python environment such as Jupyter Notebook or Spyder, and run each section sequentially. The output will be displayed in the console and in any visualizations.

# Future Work :

**Data Preprocessing and Feature Engineering**

One potential area for improvement is in the data preprocessing and feature engineering steps. Some possible future work includes:

- Exploring additional time-based features, such as day of the year or month of the year.
- Experimenting with different sensor-specific features, such as exponential moving average or peak detection.
- Testing different methods for handling missing or invalid data points, such as interpolation or removal based on a threshold.

**2. Model Architecture and Hyperparameter Tuning**

Another area for improvement is in the model architecture and hyperparameter tuning. Some potential future work includes:

- Trying different types of deep learning models, such as convolutional neural networks or hybrid models.
- Optimizing the hyperparameters further using more advanced methods, such as Bayesian optimization or grid search with cross-validation.
- Implementing more advanced training techniques, such as transfer learning or multi-task learning.

**3. Online Learning and Streaming Data**

A final area for improvement is in making the model more suitable for online learning and streaming data. Some possible future work includes:

- Modifying the existing model to be more adaptive to changes in traffic patterns over time.
- Developing a mechanism for incorporating new data into the existing model without the need for retraining from scratch.
- Exploring different methods for detecting and reacting to anomalies in real-time traffic data.

**Step-by-Step Guide to Implement**

To implement these future work steps, you can follow the following steps:

1. Decide which area of the project you want to improve, such as data preprocessing, model architecture, or online learning.
2. Research and experiment with different methods or techniques for improving that area of the project, either by modifying the existing code or writing new code.
3. Evaluate the performance of the new code using metrics such as mean squared error or precision-recall curve.
4. Iterate on the new code until you achieve the desired level of performance.
5. Integrate the new code into the existing project and rerun the entire pipeline to see the impact of the changes on the final output.

Overall, by implementing these future work steps, you can continue to improve the accuracy and robustness of the traffic flow prediction model and make it more suitable for real-world applications.

# Exercise :

**Try to answers the following questions by yourself to check your understanding for this project. If stuck, detailed answers for the questions are also provided.**

1. **What are some potential areas for improvement in the data preprocessing and feature engineering steps of this project?**
   One potential area for improvement is in exploring additional time-based features, such as day of the year or month of the year. Another area for improvement is in experimenting with different sensor-specific features, such as exponential moving average or peak detection. Additionally, different methods for handling missing or invalid data points, such as interpolation or removal based on a threshold, could also be explored.

2. **What are some potential areas for improvement in the model architecture and hyperparameter tuning steps of this project?**
   One potential area for improvement is in trying different types of deep learning models, such as convolutional neural networks or hybrid models. Additionally, optimizing the hyperparameters further using more advanced methods, such as Bayesian optimization or grid search with cross-validation, could be explored. Finally, more advanced training techniques, such as transfer learning or multi-task learning, could also be implemented.

3. **What is the purpose of the LSTM layer in the traffic flow prediction model?**
   The LSTM layer in the traffic flow prediction model is responsible for learning the temporal dependencies between traffic flow values at different time steps. This allows the model to capture long-term patterns in the traffic data and make accurate predictions.

4. **What is the purpose of evaluating the model on the testing data after training it on the training data?**
   Evaluating the model on the testing data after training it on the training data allows us to estimate the model's performance on unseen data. This helps to ensure that the model has not overfit to the training data and is able to make accurate predictions on new data.

5. **What is the purpose of visualizing the actual and predicted traffic flow values on a graph?**
   Visualizing the actual and predicted traffic flow values on a graph allows us to compare the model's predictions to the ground truth values and visually inspect the model's performance. This can help us to identify patterns and anomalies in the data that the model may have missed, and can also be useful for communicating the results to stakeholders or other team members.