

## Introduction:

This project is a part of IIT Dh's Tech fest 'PARSEC'. This project is supposed to predict the solar radiation and geomagnetic storms in the near future if past data is given.

This report delves into the intricate task of determining Disturbance Storm Index (DST) values through a comprehensive data analysis approach.

It also aims to provide the reader with all the information such as why we took some particular features to scale, preprocess data, scale it and our thoughtprocess while choosing our Model.

Our objective for the Quest 1 of AURORA is to predict Dst (Disturbance Storm-TimeIndex), which is a measure of Solar Storm Intensity, using Level 2 Data

refined from DSCOVR's Raw Data that

- Push the prediction performance envelope.
- Is within operationally feasible restrictions
- Is making use of specific real-time solar-wind data sources.

## Data Preprocessing:

1. Importing required libraries for smooth functioning of our code.

2. Reading each csv file, checking data for prediction.

->Start of preprocessing

1. Timedelta columns and indices readjusted/periods were grouped,max-min values,mean and standard deviations observed.

2. Checking for null values & Visualising the different parameters through line plots, We infer the following :

-> The proportion varied by feature, but every one required some kind of imputation[a huge number of na values encountered for each of the feature].

-> The na values were imputed using the forward fill method(ffill).

-> The features 'bx\_gse' and bx\_gsm are closely related. Their respective plots showed near resemblance.

-> The source column was dropped to find the correlation of other data.

-> A correlation matrix was plotted against all the features.

=>speed and temp. were highly anti-correlated with correlation values ranging from [-0.6,-0.4].

## Correlation And Data Visualisation:

1. Highly correlated features were shortlisted for feature scaling and next procedures...

-> "bt"

-> "temperature"

-> "bx\_gse"

-> "by\_gse"

-> "bz\_gse"

-> "speed"

-> "density"

-> "sunspots numbers"

2. Subsets of these data were made and monthly/daily data was made hourly so as to aggregate with other features

Missing values when made hourly were filled using 'mean'.

'smoothed\_ssn' was imputed using forward fill.

'solar\_wind' was imputed through interpolation.

3. Standard Scaler was used for Feature Scaling: to merged and fitted the solar\_wind and sunspots data
4. Then the correlated features were simply preprocessed using the preprocess\_features function within Standard Scaler.
5. Modifying the shape of labels dst by data shifting as we want to predict dst of 1 hour ahead also by creating variables t0 and t1 (we are basically creating two columns t0 and t1 which will predict the dst values with 1 hour gap)

Hence we can do multi-step prediction by providing it both steps.

-List of column names for the labels joining the t0 and t1 columns with the dst dataframe

6. The given data is large in size. So, the data is splitted in three different periods namely, train, test and validation

we have created a new dataframe interim which includes all the rows from the original dataframe "data" except for those whose index values matches with the index of the dataframe "test".

7. Finally, Data splits: 3 parts(width=0.75)--> "train\_a", "train\_b" and "train\_c" were created from the existing preprocessed dataset so as to monitor the train, validation and test split counts over hourly timesteps. using plots.

## MODEL TRAINING-->

1. We chose the Tensorflow neural network library for our model training:

-> Tensorflow provides functionality for working with timeseries data.

-> In this case, the code is utilizing the timeseries\_dataset\_from\_array function from tensorflow.keras.preprocessing to create time series datasets.

-[import tensorflow as tf]

-[from tensorflow.keras.preprocessing import timeseries\_dataset\_from\_array]

-> The data\_config dictionary holds configuration parameters such as the number of time steps (timesteps) and batch size (batch\_size).

-> The timeseries\_dataset\_from\_df function takes a DataFrame (df) and batch size as input and creates a time series dataset. It does this by iterating over periods in the DataFrame, extracting input and output sequences, and using timeseries\_dataset\_from\_array to create the dataset.

-> The train\_ds and val\_ds are created by calling the timeseries\_dataset\_from\_df function for the training and validation data, respectively.

-> Finally, the code prints the number of batches in the training and validation datasets.

2. We have used LSTM (Long Short-Term Memory) Neural network architecture for using Sequential Modeling. Since our task is a type of time series prediction task, this RNN is the optimal one.

-The Sequential model is created, representing a linear stack of layers.

-> We have added the Dense Layer with number of neurons equal to len(YCOLS)

-> We have configured the model after several trials for n\_epochs=16, n\_neurons=512, (an epoch refers to one complete pass through the entire training dataset during the training phase)

For the compilation of the model, we have used the Adam Algorithm and the loss is calculated by mean\_squared error.

3. Prints only the shape of the first batch.

4. The fit method trains the model on the provided training data (train\_ds) for the specified number of epochs, using the specified batch size. The training progress is monitored, and if a validation dataset is provided (val\_ds), the model's performance on the validation set is evaluated after each epoch.

5. The training history, including metrics like loss and accuracy on both the training and validation sets, is returned and stored in the history variable. This information can be useful for later analysis and visualization.

on of the training process.

## 6. Finishing ===>

-> Plotting Training History: it includes metrics like loss and other metrics specified during training, over epochs.

-> Preparing Test Dataset: A time series dataset for testing (test\_ds) is created using the timeseries\_dataset\_from\_df function, similar to what was done for the training and validation datasets.

-> Model Evaluation on Test Dataset: The trained model is evaluated on the test dataset (test\_ds) using the evaluate method. The resulting Mean Squared Error (MSE) is then printed. Additionally, the Root Mean Squared Error (RMSE) is calculated and printed, which is a common metric for regression problems.

-> Saving Model and Scaler: The trained model is saved to a file named "model". The scaler used for normalization during training (scaler) is saved to a file named "scaler.pck".

-> Updating and Saving Configuration: The data\_config dictionary is updated by adding a new key, "solar\_wind\_subset", and its value (SOLAR\_WIND\_FEATURES).

The updated data\_config is then saved to a JSON file named "config.json". This file stores configuration information that might be needed when deploying or using the model in different contexts.