**Part A**

**Question 1: Explain below the 5 components shown in orange boxes. Explain which Azure components you will use where in this big data architecture and why.**
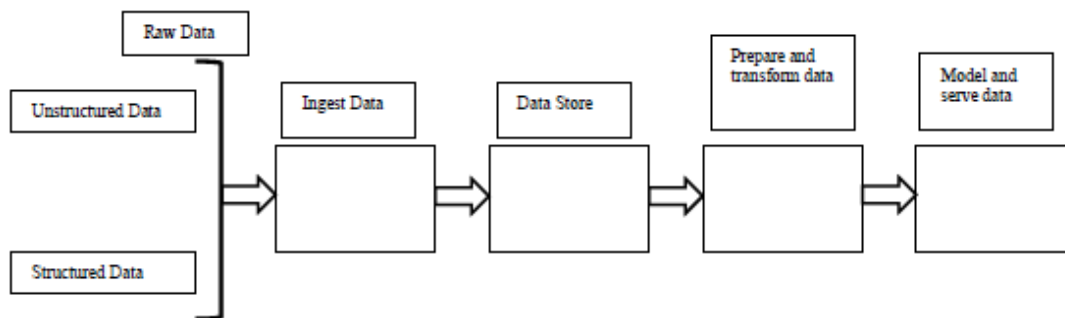
Azure Data Lake: Provides scalable storage for relational, semi-relational, and/or unstructured data

Azure Databricks: Data analytics platform optimized for Microsoft cloud services platform

Azure Data Factory: Cloud-based data integration service for Extract, Load, and Transform pipelines

Azure Synapse Analytics: Analytics engine designed to process big data quickly using parallel processing architecture

Azure Cosmos DB: Globally distributed, massively scalable, multi-model NoSQL database service for modern app development.



To ingest data: **Azure Data Factory** offers a scale-out, managed data movement solution. It allows users to ingest data with Azure Machine Learning.

Data store: **Azure Data Lake** allows the storage of any type of data of any size while Azure Cosmos DB only stores NoSQL data. Since the architecture deals with both structured and unstructured data, Azure Data Lake is used to store any ingested data.
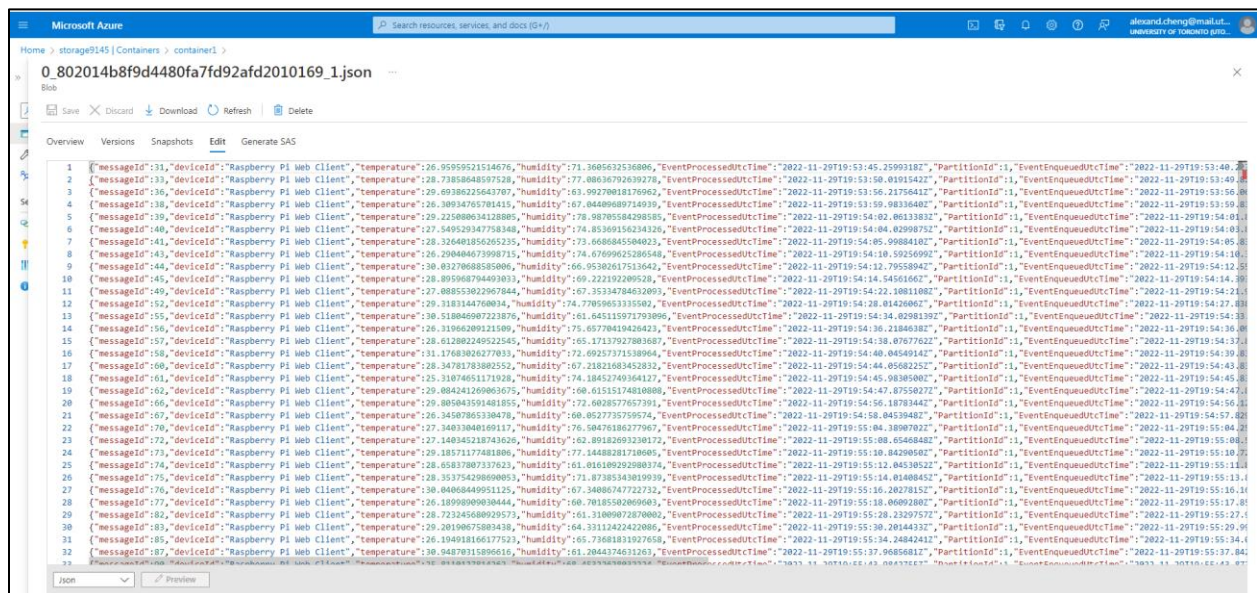
Prepare and Transform Data: **Azure Databricks** allows the preparation and the transformation (i.e. clean, sort, merge, join, etc.) of data in a notebook environment. It is used to prepare and transform the stored data.

Model and serve data: **Azure Synapse Analytics** are a limitless analytics service able to fulfill the needs of complex big data analysis and modeling.
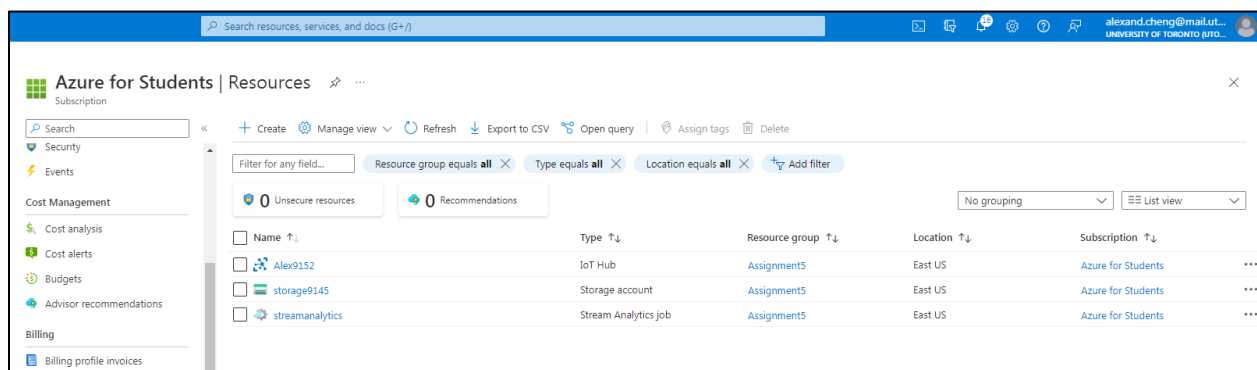
**Question 2: Explain how Stream Analytics works in Azure.**

Azure stream analytics can be used for real-time data ingestion and processing in the Azure cloud. Azure Stream Analytics allows the user to input (ingest) real time data, query (analyze), and an output (deliver). Data is ingested from Azure Event Hubs, IoT Hub, or Blob Storage. It allows simple transformation and queries to be applied on the data as well as machine learning. Analysis can be performed based on SQL query language for filtering, sorting, aggregating, and joining streaming data over a period of time. Transformed data can be sent to other services such as Azure Function, Service Bus, or event hubs, or can be sent to Power BI for visualization purposes and real-time dashboarding. Users can also store the data in other Azure storage services for training machine learning models.
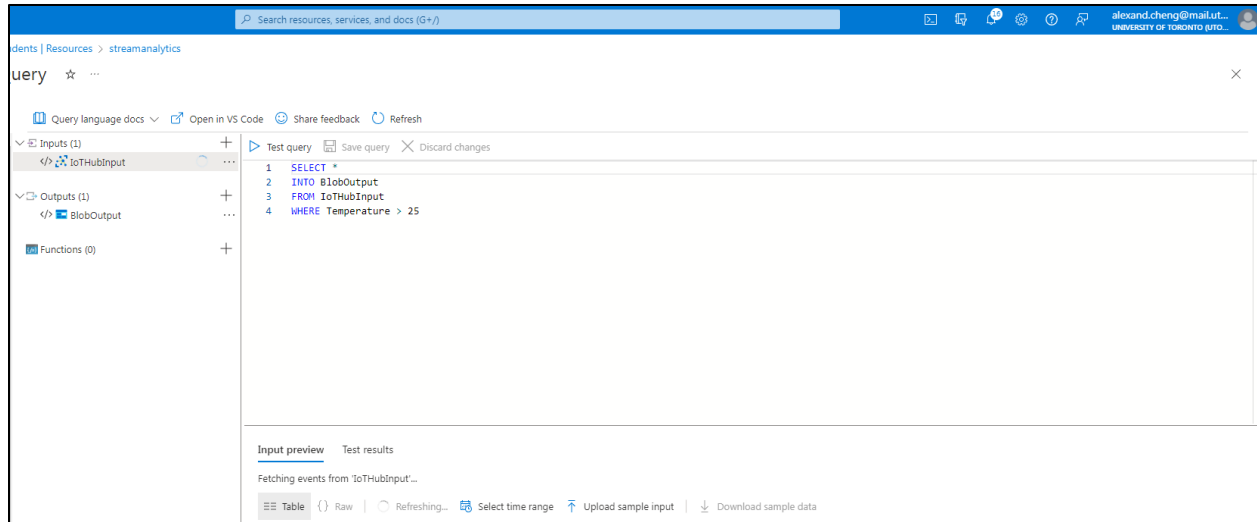
**Question 3: Implementing a Stream Analytics job by using the Azure portal**
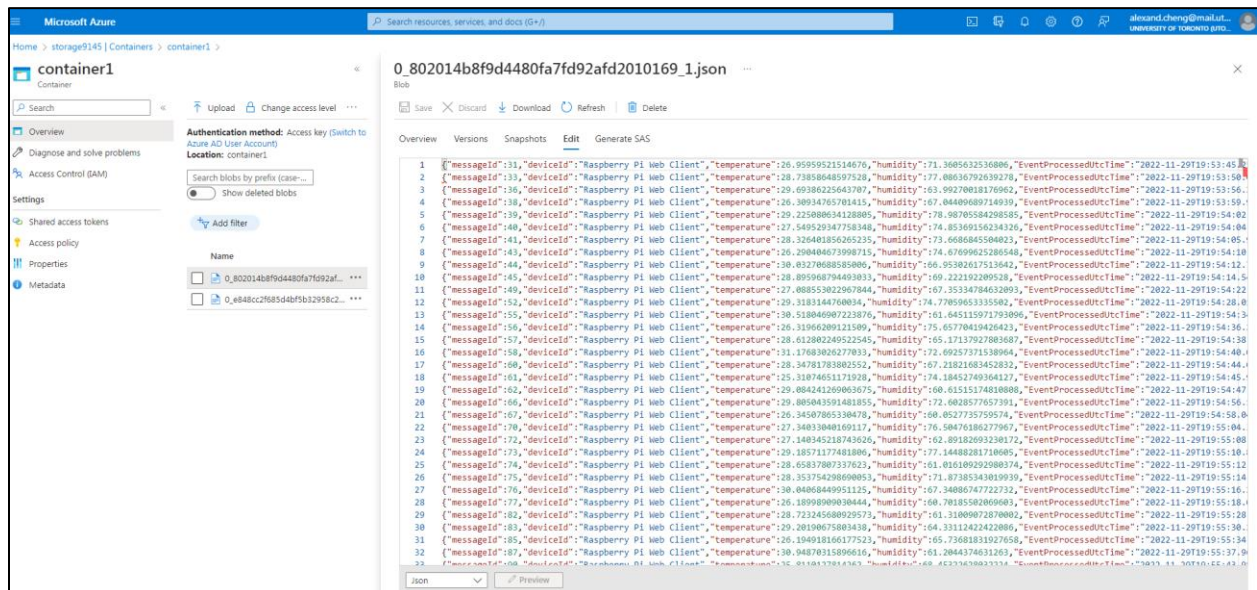


Screenshot of the top 30 results



Deployed Resources

Screenshot of Query



Another screenshot of the top 30 results

**Part B – Implement in Azure Machine Learning using Azure Notebook**

https://archive.ics.uci.edu/ml/datasets/Gas+Turbine+CO+and+NOx+Emission+Data+Set

1. Explain what problem you are going to solve using this dataset. Provide a brief overview of your problem statement. [Discuss your problem statement with your TAs if they approve then you can proceed with the next steps.]

**Problem Statement:**

The gas turbine is very useful in the modern world since it is relatively compact in size and creates a lot of power. Gas turbines are used in backup power systems in Manhattan, for example, when the grid goes down due to natural disasters, gas turbines power up and can produce power for emergencies.

The dataset contains 36733 instances of 11 sensor measures aggregated over one hour (by means of average or sum) from a gas turbine located in Turkey's northwestern region for the purpose of studying flue gas emissions, namely CO and NOx (NO + NO2). Some of the features listed are ambient temperature, pressure, humidity, etc. We can examine which main features are used for predicting hourly net energy yield.

By determining which features are the most significant in determining the Turbine energy yield, we can determine if the environment is suitable for the energy yield required (clients will have a minimum energy output needed if used as a backup), this can save a lot of cost by researching the environment and critical components before implementing and finding out that the location chosen is not suitable. By comparing experimental results with theoretical results for gas turbine energy yield, the user can determine if there's any external factors that are significant that are not being included.
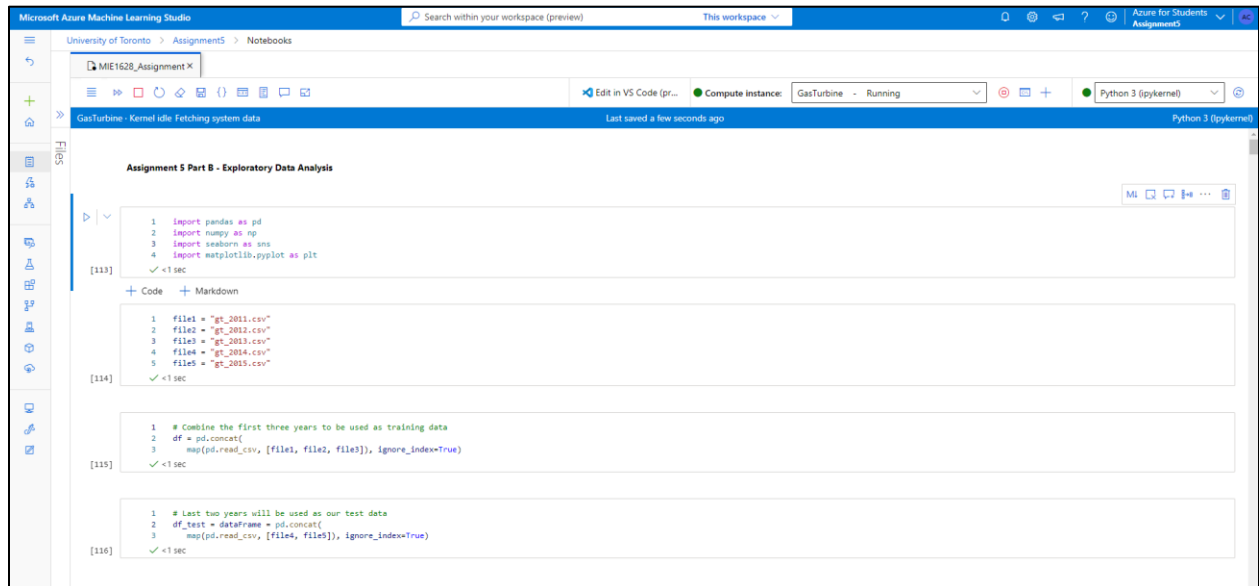
There will be no data cleaning required since there are no N/A or missing values inside the dataset, however, the data is separated by years throughout 5 different years. We will merge the first three years to be used as our training data and the last two years will be used as our test data to determine how accurate our model is, this will be a 60/40 split. Since the data does not require data cleaning, more effort will be placed on exploratory data analysis and modelling. We will spend a portion on examining the gas emissions (CO and NOx), due to the rising concerns of climate change and renewable energy.

2. Explain your dataset. Explore your dataset and provide at least 5 meaningful charts/graphs with an explanation.
3. Do data cleaning/pre-processing as required and explain what you have done for your dataset and why?

A summary of the Exploratory Data Analysis and screenshots will be shown: Please look at the ipynb file for a detailed explanation of each step, otherwise there will be too many screenshots.

**List of Exploratory Data Analysis for the Dataset**

1. Correlation Plot - Can see which features are related with one another
2. Histogram - see the count distribution of each feature
3. Plot each feature against TEY to see the direct relationship
4. Pairplot to see relationship between all features
5. Time Series to compare changes/differences between years for TEY and CO (Climate Change) to see significance differences



Deployed compute instance to run notebook on Azure Portal

**From the dataset webpage - Meaning of each acronym**

Variable (Abbr.) Unit Min Max Mean

Ambient temperature (AT) C 6.23 37.10 17.71

Ambient pressure (AP) mbar 985.85 1036.56 1013.07

Ambient humidity (AH) (%) 24.08 100.20 77.87

Air filter difference pressure (AFDP) mbar 2.09 7.61 3.93

Gas turbine exhaust pressure (GTEP) mbar 17.70 40.72 25.56

Turbine inlet temperature (TIT) C 1000.85 1100.89 1081.43

Turbine after temperature (TAT) C 511.04 550.61 546.16

Compressor discharge pressure (CDP) mbar 9.85 15.16 12.06

Turbine energy yield (TEY) MWH 100.02 179.50 133.51

Carbon monoxide (CO) mg/m3 0.00 44.10 2.37

Nitrogen oxides (NOx) mg/m3 25.90 119.91 65.29

The dataset contains instances of 11 sensor measures aggregated over one hour (by means of average or sum) throughout 5 years

---

1    df

✓ < 1 sec

| | AT | AP | AH | AFDP | GTEP | TIT | TAT | TEY | CDP | CO | NOX |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 4.5878 | 1018.7 | 83.675 | 3.5758 | 23.979 | 1086.2 | 549.83 | 134.67 | 11.898 | 0.32663 | 81.952 |
| 1 | 4.2932 | 1018.3 | 84.235 | 3.5709 | 23.951 | 1086.1 | 550.05 | 134.67 | 11.892 | 0.44784 | 82.377 |
| 2 | 3.9045 | 1018.4 | 84.858 | 3.5828 | 23.990 | 1086.5 | 550.19 | 135.10 | 12.042 | 0.45144 | 83.776 |
| 3 | 3.7436 | 1018.3 | 85.434 | 3.5808 | 23.911 | 1086.5 | 550.17 | 135.03 | 11.990 | 0.23107 | 82.505 |
| 4 | 3.7516 | 1017.8 | 85.182 | 3.5781 | 23.917 | 1085.9 | 550.00 | 134.67 | 11.910 | 0.26747 | 82.028 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 22186 | 4.8631 | 1027.0 | 81.084 | 4.2825 | 34.045 | 1100.0 | 529.98 | 168.38 | 14.290 | 1.25380 | 78.397 |
| 22187 | 4.5173 | 1027.4 | 80.813 | 4.2481 | 33.904 | 1100.1 | 530.47 | 168.07 | 14.344 | 1.08080 | 78.251 |
| 22188 | 4.2717 | 1027.9 | 80.380 | 4.2817 | 34.165 | 1099.9 | 529.56 | 168.55 | 14.395 | 1.04720 | 77.269 |
| 22189 | 4.0853 | 1028.6 | 78.907 | 4.2313 | 33.802 | 1100.1 | 530.61 | 167.98 | 14.343 | 1.08750 | 77.985 |
| 22190 | 4.2148 | 1029.4 | 70.679 | 4.2049 | 33.768 | 1100.0 | 530.97 | 167.30 | 14.291 | 1.13370 | 78.950 |

22191 rows × 11 columns

Training Data, the target variable is the Turbine Energy Yield (TEY)

```
    1    # Basic information
    2    df.info()
    ✓ <1 sec
```
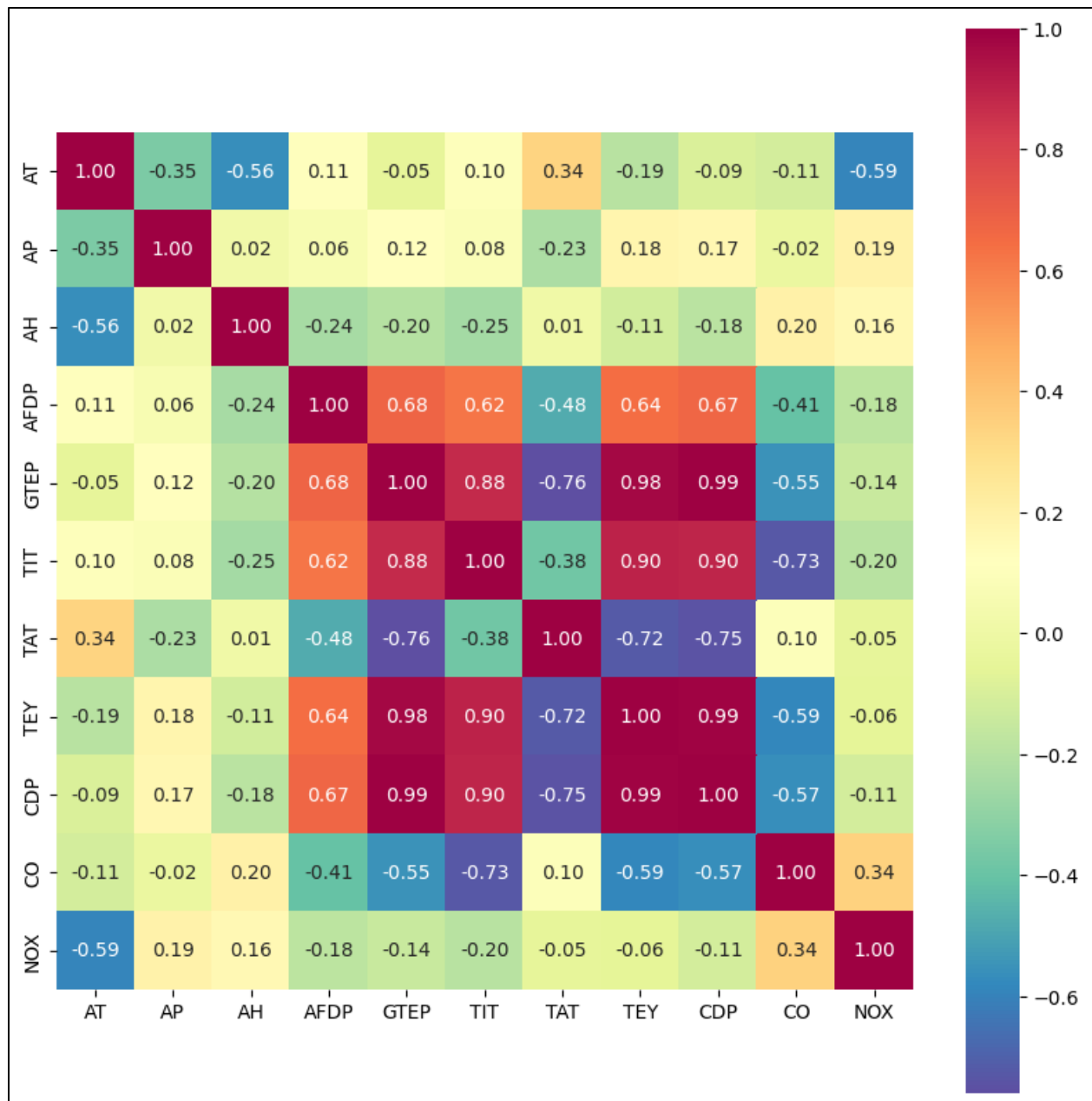
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 22191 entries, 0 to 22190
Data columns (total 11 columns):
 #    Column   Non-Null Count    Dtype
---   ------   --------------    -----
 0    AT       22191 non-null    float64
 1    AP       22191 non-null    float64
 2    AH       22191 non-null    float64
 3    AFDP     22191 non-null    float64
 4    GTEP     22191 non-null    float64
 5    TIT      22191 non-null    float64
 6    TAT      22191 non-null    float64
 7    TEY      22191 non-null    float64
 8    CDP      22191 non-null    float64
 9    CO       22191 non-null    float64
 10   NOX      22191 non-null    float64
dtypes: float64(11)
memory usage: 1.9 MB
```

We can observe that our data does not require any data cleaning/pre-processing

**Correlation Plot**



We can observe that multiple features are correlated with one another, mostly the compressor features and not regarding the environment conditions.

Features such as the air filter difference pressure, gas turbine exhaust pressure, turbine inlet temperature, turbine energy yield, and compressor discharge pressure are highly correlated with one another.

For the Turbine energy yield, we can observe that the values are highly correlated with the gas turbine exhaust pressure, turbine inlet temperature, and compressor discharge pressure.

**Histogram**



We can observe that there are some features that have a similar distribution to TEY such as the Compressor Discharge Pressure (CDP) and a bit of Gas Turbine Exhaust Pressure (GTEP), we can also observe that Ambient Pressure (AP) and NOx follow like a normal distribution based on the bell-shaped curve

**Probability Plots**





We can observe that the ambient pressure there follows a normal distribution, the amount of NOx emissions has a normal distribution around the average values and skews around the ends.

**Find the relationship between features and Turbine Energy Yield using Scatter Plots**



We can double check the relationship by plotting it against TEY, it is observed that GTEP and CDP fall in a straight line, thus follow a similar distribution as the TEY. Thus, we can be confident that GTEP and CDP will be key features on predicting the TEY.

**Gas Emissions and Global Warming Concerns**





There's also an objective to minimize the emissions of gas turbines while maintaining the purpose of the gas turbine, high yield energy. It seems that a high energy yield does not mean higher CO or NOx emissions

**Pairplots –** We're looking for graphs that fit a linear relationship with one another



We can see a linear relationship between GTEP and TEY and CDP, so we know those features (GTEP and CDP) will have more feature importance in our machine learning models to predict the TEY.

**Time Series** – See if there's any difference between years 2011, 2012, and 2013 since we'll be merging them into 1 dataset as our training data. Each index roughly represents the measurements aggregated over one hour (by means of average or sum) from a gas turbine. Here we compare the TEY over the three years.

We can see that the minimum and maximum range for the Turbine Energy Yield stays the same throughout the years, we see that the average energy output is slightly decreasing. In 2011, we can observe that the Turbine energy yield was higher for the first quarter, and for the middle of the graphs, there are many more low points, this may be since we require more energy/electricity in the winter compared to the summer.

Compare CO emissions over the three years

We can see that the minimum and maximum range for the CO emissions are roughly the same throughout the years, we see that the average CO emissions is increasing over the years. There are a few outliers with high spikes in NOx emissions.

4. Implement 2 machine learning models and explain which algorithms you have selected and why. Compare them and show success metrics (Accuracy/RMSE/Confusion Matrix) as per your problem. Explain results.

**Part 4 - Machine Learning Models**

```
1    #Standardize Data
2    mu,si = df.mean(), df.std() #Calculate the overall mean and standard deviation of the quality scores
3    df_standard = df - mu #Subtract the mean from every entry
4    df_standard = df_standard / si #Divide every entry by the standard deviation
5    #Confirm if the data has been normalized
6    df_standard
```
✓ <1 sec

|        | AT        | AP       | AH        | AFDP      | GTEP      | TIT      | TAT       | TEY      | CDP       | CO        | NOX      |
|--------|-----------|----------|-----------|-----------|-----------|----------|-----------|----------|-----------|-----------|----------|
| 0      | -1.784962 | 0.920396 | 0.296067  | -0.570152 | -0.316208 | 0.185248 | 0.559086  | 0.070673 | -0.145574 | -0.822286 | 1.193838 |
| 1      | -1.825028 | 0.857863 | 0.336311  | -0.576199 | -0.322821 | 0.179310 | 0.587625  | 0.070673 | -0.150958 | -0.769488 | 1.232347 |
| 2      | -1.877892 | 0.873496 | 0.381083  | -0.561512 | -0.313610 | 0.203062 | 0.605786  | 0.097504 | -0.016341 | -0.767920 | 1.359112 |
| 3      | -1.899775 | 0.857863 | 0.422477  | -0.563981 | -0.332268 | 0.203062 | 0.603192  | 0.093136 | -0.063008 | -0.863911 | 1.243946 |
| 4      | -1.898687 | 0.779696 | 0.404367  | -0.567313 | -0.330851 | 0.167434 | 0.581139  | 0.070673 | -0.134804 | -0.848055 | 1.200724 |
| ...    | ...       | ...      | ...       | ...       | ...       | ...      | ...       | ...      | ...       | ...       | ...      |
| 22186  | -1.747520 | 2.217963 | 0.109865  | 0.302077  | 2.061130  | 1.004688 | -2.015924 | 2.174115 | 2.001133  | -0.418422 | 0.871717 |
| 22187  | -1.794550 | 2.280496 | 0.090390  | 0.259620  | 2.027829  | 1.010626 | -1.952360 | 2.154772 | 2.049596  | -0.493778 | 0.858488 |
| 22188  | -1.827952 | 2.358663 | 0.059272  | 0.301090  | 2.089471  | 0.998750 | -2.070408 | 2.184723 | 2.095366  | -0.508414 | 0.769508 |
| 22189  | -1.853303 | 2.468096 | -0.046585 | 0.238885  | 2.003739  | 1.010626 | -1.934198 | 2.149156 | 2.048698  | -0.490860 | 0.834385 |
| 22190  | -1.835691 | 2.593163 | -0.637888 | 0.206301  | 1.995709  | 1.004688 | -1.887498 | 2.106725 | 2.002031  | -0.470736 | 0.921825 |

22191 rows × 11 columns

Standardized Training data and remove the target column (TEY)

**Multiple Linear Regression Model**

Multiple Linear Regression model was chosen to predict the Turbine energy yield, this is one of the simplest predictive regression models. We can see if a linear scaling of the features can be summed up to determine the turbine yield energy. Linear regression predicts the weights of each feature. It used the method of ordinary least squares, minimize the sum of squared residuals (SSR) to determine the best weights.

```
1    from sklearn.linear_model import LinearRegression
2    model = LinearRegression().fit(df_standard, y_train)
   ✓ <1 sec
```

```
1    y_prediction =  model.predict(df_test_standard)
2    from sklearn.metrics import r2_score
3    from sklearn.metrics import mean_squared_error
4    # predicting the accuracy score
5    score=r2_score(y_test,y_prediction)
6    print("r2 score is ",score)
7    print("mean_sqrd_error is==",mean_squared_error(y_test,y_prediction))
8    print("root_mean_squared error of is==",np.sqrt(mean_squared_error(y_test,y_prediction)))
   ✓ <1 sec
```

```
r2 score is  0.9831772531448139
mean_sqrd_error is== 3.7727061796151267
root_mean_squared error of is== 1.9423455355870969
```

The accuracy score is high at about 98.31%, the RMSE is also low at about 1.94 MWH considering its goes on average to 130 MWH. R2 score is the coefficient of determination, the amount of variation in the target (Turbine energy yield) that can be explained by the dependence of the input features. We can further improve the accuracy by hyper tuning the model. But this would be easier to leave the hyper parameter tuning to AutoML in Part B-5.

```
1    y_check = model.predict(df_standard)
2    score=r2_score(y_check,y_train)
3    print("r2 score is ",score)
4    print("mean_sqrd_error is==",mean_squared_error(y_check,y_train))
5    print("root_mean_squared error of is==",np.sqrt(mean_squared_error(y_check,y_train)))
   ✓ <1 sec
```
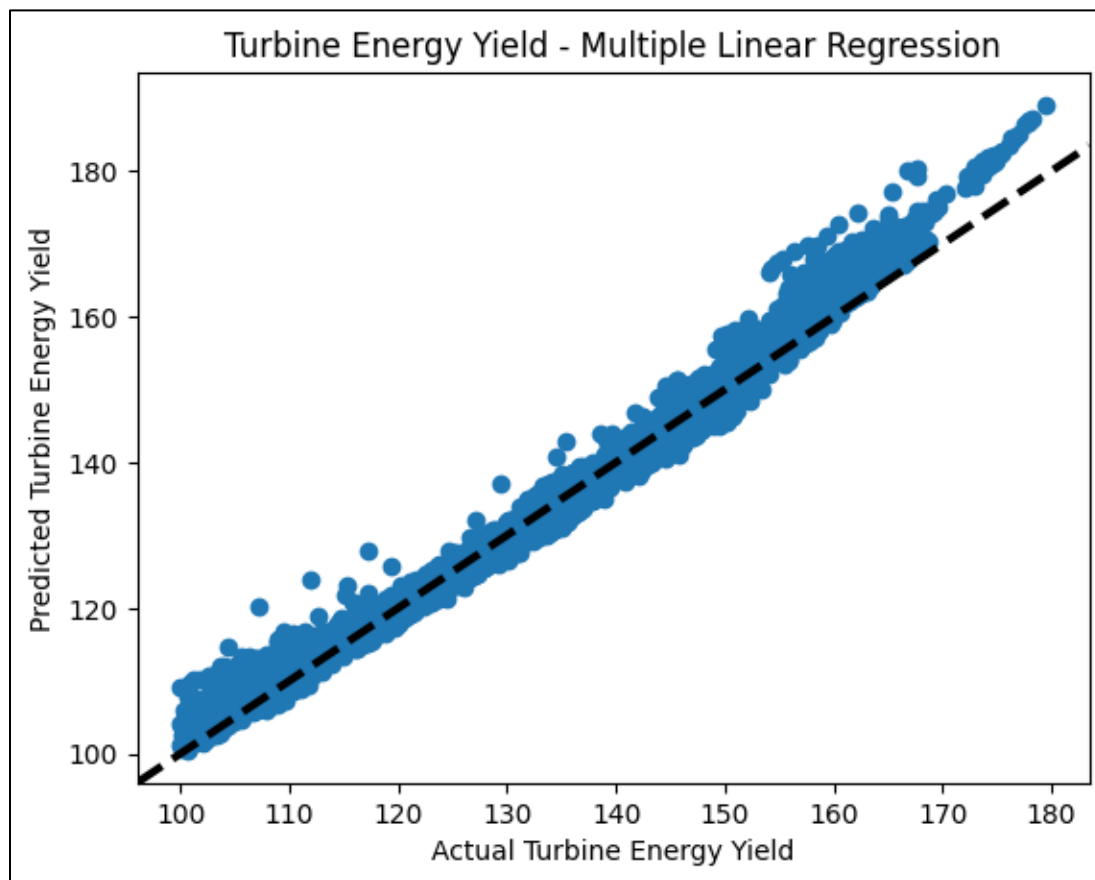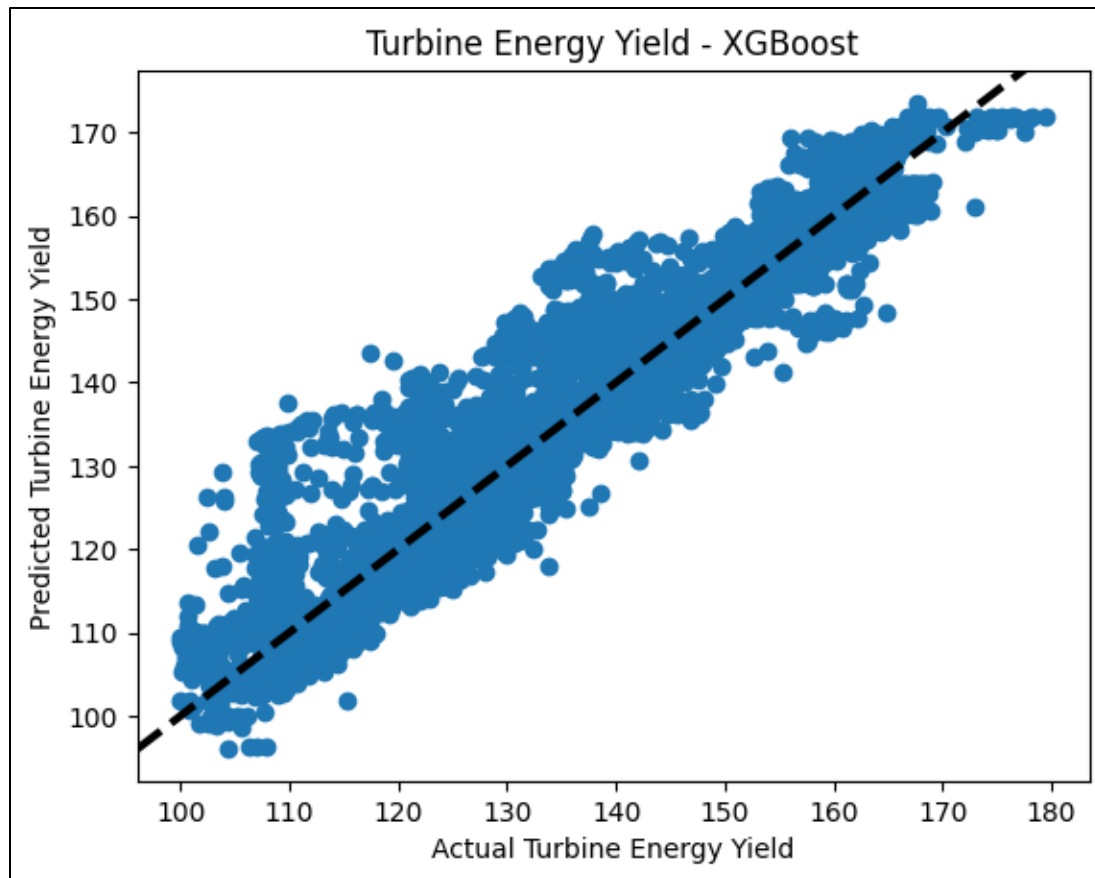
```
r2 score is  0.9978385437961127
mean_sqrd_error is== 0.5539177168382373
root_mean_squared error of is== 0.7442564859228553
```

We ran the model against the training data again to see whether it was overfitting or underfitting, the model did slightly better on the training data compared to the test data so the model but still has relatively high R2 scores. The simple multi linear regression model is slightly overfitting the training data.

```
1    print(model.coef_)
2    print(df_standard.columns)
3    print(model.intercept_)
   ✓ <1 sec
```

```
[-2.60817789 -0.3806715  -0.0876376  -0.06562801  0.93759214  9.64556195
 -4.29575866  2.88959302 -0.03448007 -0.20546109]
Index(['AT', 'AP', 'AH', 'AFDP', 'GTEP', 'TIT', 'TAT', 'CDP', 'CO', 'NOX'], dtype='object')
133.5373930872805
```

Since the data is normalized, the bias is set at around the average value. We can examine the weight of each feature. We can observe that a higher Turbine inlet temperature and Compressor Discharge pressure is correlated for a higher turbine energy yield. We can also see that Turbine after temperature is negatively correlated with a higher energy yield (most of the energy turning into heat possibly). \

**XGBoost**

XGBoost stands for Extreme Gradient Boost, boosting is an ensemble technique where new models are added to correct the errors made by existing models. Models are added sequentially until no further improvements can be made. Gradient boosting is an approach where new models are created that predict the residuals or errors of prior models and then added together to make the final prediction. It is called gradient boosting because it uses a gradient descent algorithm to minimize the loss when adding new models. XGBoost is a library for developing fast and high-performance gradient boosting tree models which can be used for solving regression type problems.

```
1    import xgboost as xgb
```
✓ <1 sec

```
1    xg_reg = xgb.XGBRegressor(objective ='reg:squarederror', colsample_bytree = 0.3, learning_rate = 1,
2                              max_depth = 5, alpha = 10, n_estimators = 10)
```
✓ <1 sec

\+ Code    \+ Markdown

```
1    xg_reg.fit(df_standard,y_train)
2
3    preds = xg_reg.predict(df_test_standard)
```
✓ <1 sec

```
1    score=r2_score(y_test,preds)
2    print("r2 score is ",score)
3    print("mean_sqrd_error is==",mean_squared_error(y_test,preds))
4    print("root_mean_squared error of is==",np.sqrt(mean_squared_error(y_test,preds)))
```
✓ <1 sec

```
r2 score is  0.9138600189298023
mean_sqrd_error is== 19.31794145706231
root_mean_squared error of is== 4.395218021561878
```

First tried this model with the learning rate as 0.1, and got the RMSE at around 46, this meant that the model was overfitting. Now the RMSE is around 4.395 MWH which is relatively low with a R2 score of 91.3% without too many hyperparameter tuning.

```
1    y_checkXG = xg_reg.predict(df_standard)
2    score=r2_score(y_checkXG,y_train)
3    print("r2 score is ",score)
4    print("mean_sqrd_error is==",mean_squared_error(y_checkXG,y_train))
5    print("root_mean_squared error of is==",np.sqrt(mean_squared_error(y_checkXG,y_train)))
```
✓ <1 sec

```
r2 score is  0.9939731041110403
mean_sqrd_error is== 1.5545162976502156
root_mean_squared error of is== 1.246802429276674
```

We ran the model against the training data again to see whether it was overfitting or underfitting, the model did slightly better on the training data compared to the test data so the model but still has high R2 scores. The more complex XGBoost model is slightly overfitting the training data.

**Comparing the two-machine learning models**

Turbine Energy Yield - XGBoost

Our simple multiple linear regression model performed better than the complex XGBoost model in terms of having a higher accuracy (R2 score) of (98.31% vs 91.38%) and a lower RMSE value (1.94 vs 4.395 MWH). From the plot, we can see that the multi linear regression model does a good job predicting the values except when it gets to the extremes (100-110 or 160-170 MWH). From the XGBoost model, we can see that we are overfitting our model to the training data since it has such a large spread and relatively low R2 score. If we were to hyper tune the models and run cross fold validations with the training data, the model would give better results. This will be performed in part 5 by the AutoML program.

5. Deploy a run-time pipeline for your dataset using Azure Designer Studio.
   Or do hyperparameter tuning for your algorithms. Explain your results.
   Or Use Automated ML for your data set. Explain the best model results.

**Used Automated ML for the dataset:**



AutoML Job for the Gas Turbine Dataset

List of models that Azure tested to determine which model provides the best results:
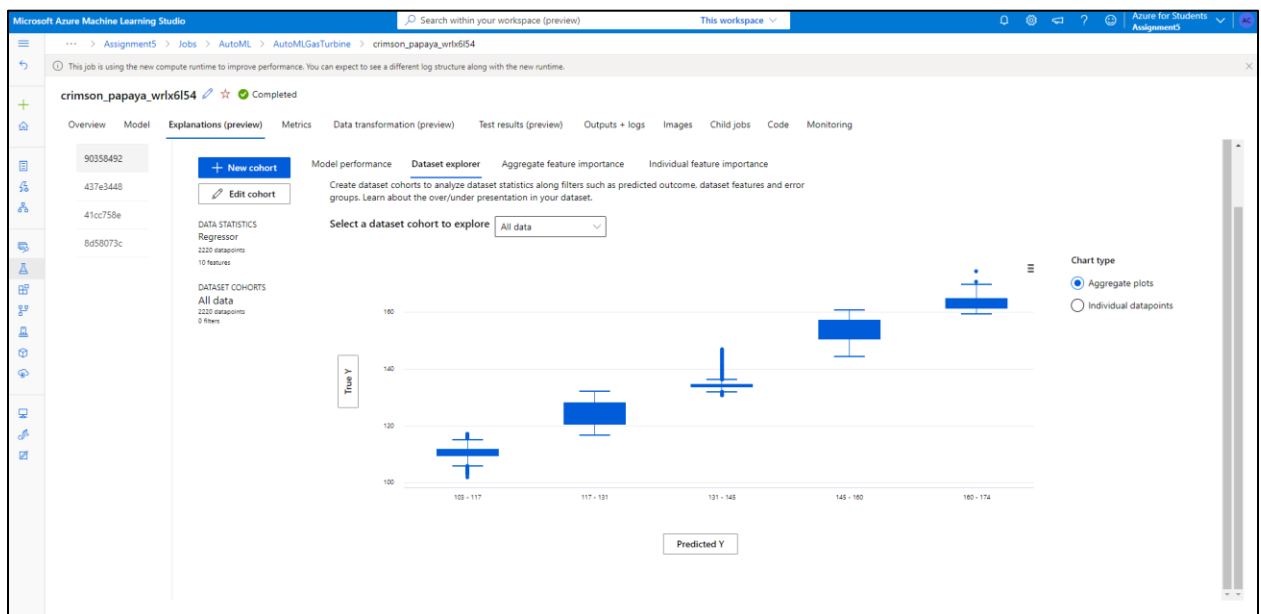
Out of the multiple models that were tested, Voting Ensemble was the best model chosen. A voting ensemble (or a "majority voting ensemble") is an ensemble machine learning model that combines the predictions from multiple other models. It is a technique that may be used to improve model performance, ideally achieving better performance than any single model used in the ensemble. We can further investigate the best model.
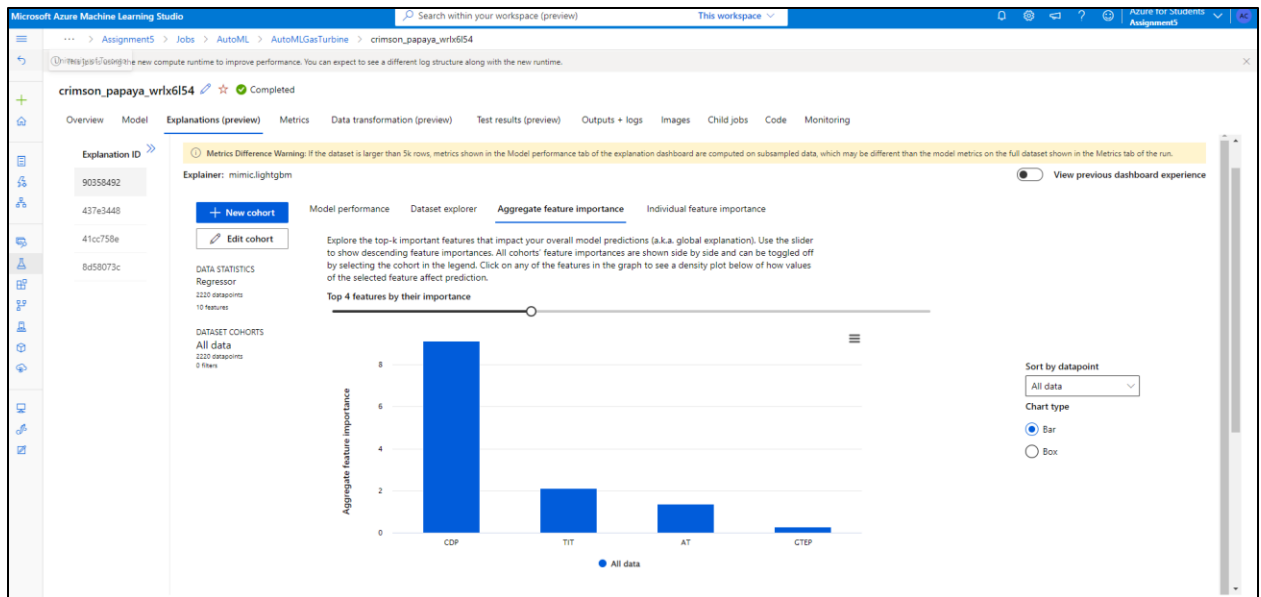
It has a very high accuracy (R2 score) and relatively small RMSE value. It can explain most the variance with the features provided.
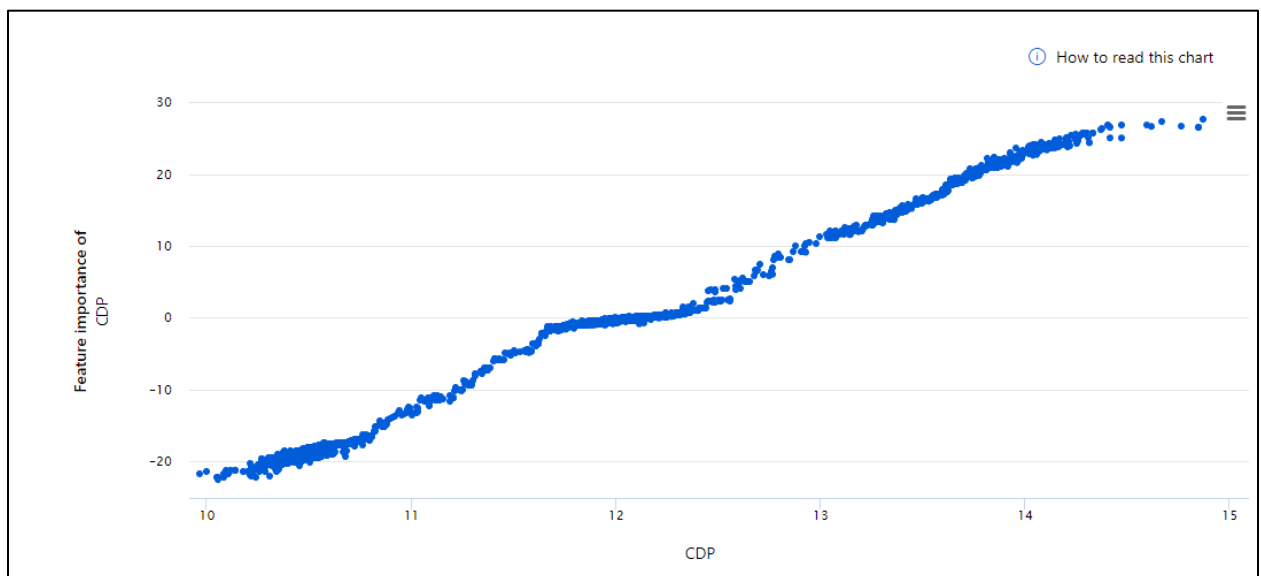
We can see that most of the mean squared error is close to 0.276, RMSE is 0.5256 MWH, and the average prediction of the energy yield is 133 MWH. There are still a few outliers predicted by the model that results in large errors. The R2 score is much higher than the multi linear regression model and the XGboost model I used in Part B-4. The RMSE is also significant smaller as well compares to those two models.
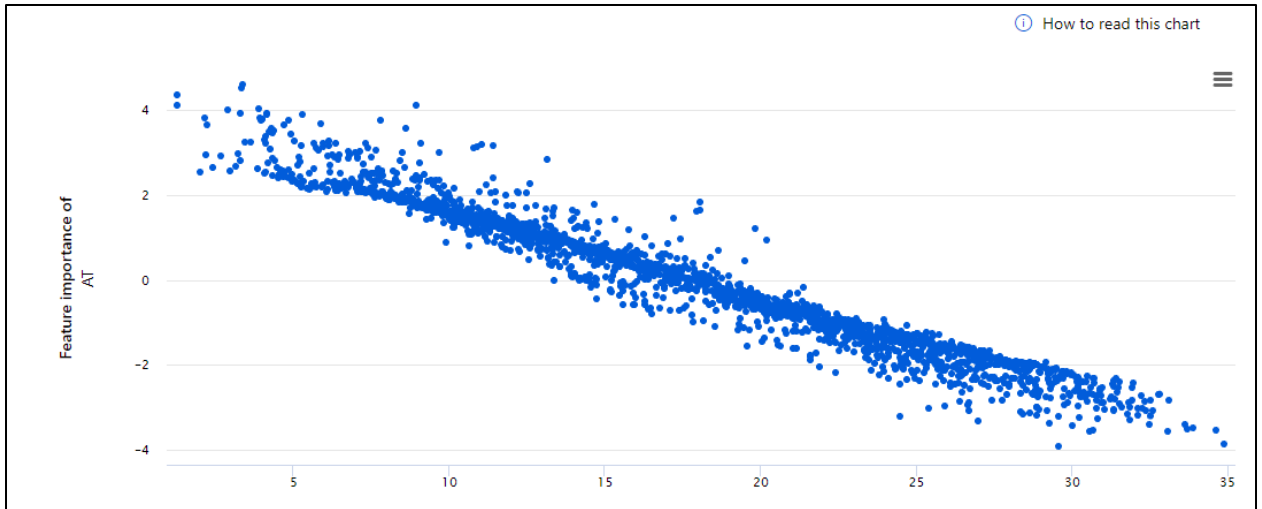


Comparing the true and predicted Turbine energy yield, we can observe that the model has the most outliers near the average value, slightly overestimating the energy yield.

We can determine the 4 most important features on determining the Turbine Energy Yield, which are: Compressor discharge pressure (CDP), Turbine inlet temperature (TIT), Ambient temperature (AT), and Gas turbine exhaust pressure (GTEP). We can see that CO and NOx emissions aren't key features on calculating the Turbine Energy Yield, thus, the next steps would be on reducing the emissions without indirectly affecting the energy yield.



We can see the feature importance of the compressor discharge pressure; a higher discharge pressure will have more significance on determining the turbine energy yield. We see that this has the largest significance on the energy yield and should be prioritized when creating/maintaining a gas turbine.

We can also see the feature importance of the ambient temperature; a lower ambient temperature will have more significance on determining the turbine energy yield.