

1. Combined Cycle Power Plant Data Set

The dataset contains data points collected from a Combined Cycle Power Plant over 6 years (2006-2011), when the power plant was set to work with full load. Features consist of hourly average ambient variables Temperature (T), Ambient Pressure (AP), Relative Humidity (RH) and Exhaust Vacuum (V) to predict the net hourly electrical energy output (EP) of the plant.

- (a) Download the Combined Cycle Power Plant data¹ from: (5 pts)
<https://archive.ics.uci.edu/ml/datasets/Combined+Cycle+Power+Plant>
- (b) Exploring the data:
 - i. How many rows are in this data set? How many columns? What do the rows and columns represent? (4 pts)
 - ii. Make pairwise scatterplots of all the variables in the data set including the predictors (independent variables) with the dependent variable. Describe your findings. (5 pts)
 - iii. What are the mean, the median, range, first and third quartiles, and interquartile ranges of each of the variables in the dataset? Summarize them in a table. (6 pts)
- (c) For each predictor, fit a simple linear regression model to predict the response. Describe your results. In which of the models is there a statistically significant association between the predictor and the response? Create some plots to back up your assertions. Are there any outliers that you would like to remove from your data for each of these regression tasks? (20 pts)
- (d) Fit a multiple regression model to predict the response using all of the predictors. Describe your results. For which predictors can we reject the null hypothesis $H_0 : \beta_j = 0$? (10 pts)
- (e) How do your results from 1c compare to your results from 1d? Create a plot displaying the univariate regression coefficients from 1c on the x-axis, and the multiple regression coefficients from 1d on the y-axis. That is, each predictor is displayed as a single point in the plot. Its coefficient in a simple linear regression model is shown on the x-axis, and its coefficient estimate in the multiple linear regression model is shown on the y-axis. (10 pts)
- (f) Is there evidence of nonlinear association between any of the predictors and the response? To answer this question, for each predictor X , fit a model of the form² (15 pts)

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \epsilon$$

¹There are five sheets in the data. All of them are shuffled versions of the same dataset. Work with Sheet 1.

²<https://scikit-learn.org/stable/modules/preprocessing.htm\#generating-polynomial-features>

- (g) Is there evidence of association of interactions of predictors with the response? To answer this question, run a full linear regression model with all pairwise interaction terms and state whether any interaction terms are statistically significant. (10 pts)
- (h) KNN Regression:
 - i. Perform k -nearest neighbor regression for this dataset using both normalized and raw features. Find the value of $k \in \{1, 2, \dots, 100\}$ that gives you the best fit. Plot the train and test errors in terms of $1/k$. (15 pts)
- (i) Can you improve your model using possible interaction terms or nonlinear associations between the predictors and response? Train the regression model on a randomly selected 70% subset of the data with all predictors. Also, run a regression model involving all possible interaction terms $X_i X_j$ as well as quadratic nonlinearities X_i^2 , and remove insignificant variables using p-values (be careful about interaction terms).³ Test both models on the remaining points and report your train and test MSEs. (Extra Credit: 15 pts)
- (j) Compare the results of KNN Regression with the linear regression model that has the smallest test error and provide your analysis. (Extra Credit: 5 pts)

³See <https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.PolynomialFeatures.html>