

<div><div>Description and motivation for the problem:</div><div>- Examine the relationship of Temperature (T), Ambient Pressure (AP), Relative Humidity (RH) and Exhaust Vacuum (V) to predict the net hourly electrical energy output (PE) of the plant. Use Random Forest and Linear Regression to predict the dependent variable PE, then compare the performance of the two algorithms.</div></div>		Table 1: Correlation Coefficient Matrix						<div><div></div><div>1</div><div>0</div><div>-1</div></div>	Table 2: Summary of Basic Statistics					
		AT		V		AP			RH		PE			
		AT		V		AP			RH		PE			
		V		AP		RH			PE		PE			
		AP		RH		PE			PE		PE			
		RH		PE		PE			PE		PE			
PE		PE		PE		PE		PE						
Initial Analysis including basic statistics:														
<div><div>- Dataset: Combined Cycle Power Plant from UCI</div><div>- The dataset has 4 continuous predictors and 1 response variable.</div><div>- No duplicate or redundant indicators were identified.</div><div>- The correlation coefficient matrix shows the potential linear relationship between the indicators (Table 1)</div><div>- Normalisation is considered as the variables are not on the same scale.</div><div>- The ambient temperature (AT) and the vacuum variable (V) have a highest chance of a linear relationship (Table 3)</div><div>- Histograms of the variables shows a decent normal distribution shape.</div><div>- Visual inspection shows insignificant outliers</div></div>														
		Table 3: Strong linear patterns												
		AT			V			AP			RH			

LINEAR REGRESSION ("LR")	RANDOM FOREST ("RF")
--------------------------	----------------------

RANDOM FOREST ("RF")

General:

- Non-parametric model
- Builds an ensemble of decision trees, each perform on a subset of features, takes an average value of the predictions - gives a numerical result.
- Can choose an attribute as a random root node or pick attribute depending on information gain.
- Considered state-of-the-art

Pros

- Allows the maintenance of high dimensionality in the dataset.
- Well regarded in terms of accuracy - particularly important in computer vision.
- Captures the non-linear nature of the data and assimilates it [2].
- Can predict both numerical and non-numerical output.
- Address overfitting problems with a reasonable number of trees.
- Handles noise and missing data well.

Cons:

- Interpretability is traded for accuracy.
- RF cannot train outside values in the range of the training data. As the predictions in random forests are provided by taking an average of the results of several trees. [4]
- Needs a large number of trees to extend its full predicting potential. However as the model increases in complexity so too does its computational cost.

Description of Choice of Training and Evaluation Methodology

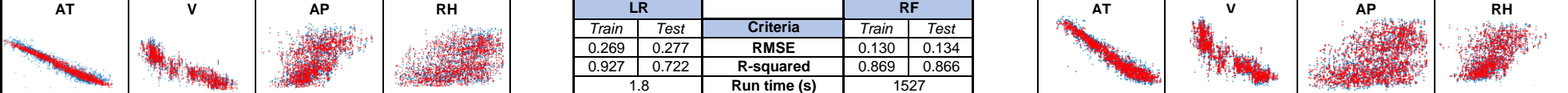
- Dataset is comprised of 9815 data points. Hold out ratio is 0.2
- Data were normalised to run the LR. In both cases the data points were initially randomised.
- Outliers are negligible so robust methods were not necessary.
- To find the optimal RF model, we train the model with the number of trees from 1 to 100 and the number of variables selected as sample from 1 to 3, which is 300 times in total.
- Models will be compared using R-squared statistics, RMSE and run-time.

LINEAR REGRESSION	RANDOM FOREST
-------------------	---------------

RANDOM FOREST

Model	Feature removal					Hold out ratio		
	N/A	AP	RH		AP, RH	0.3	0.2 - Original	0.1
LR	0.277	0.295	0.295		0.3	0.276	0.276	0.2756
RF	0.134	0.145	0.144		0.16	0.137	0.133	0.1333

Table 6: RF - Main result, predicted (red) vs. true value (blue)



Lessons learned and future work

Lesson learned:

- Both models neglect that changes may occur in the levels of interaction between a system of variables [1]. We know that PE is dependent on variables AT, V, AP & RH, but can also be certain that there are other interactions occurring within the system. These complex interactions between variables are something that are not considered in either model. For example the temperature variable (AT) can affect the interaction between other variables differently in summer and winter. Kaya et al. outlines that in this case models with localised observations can prove to be more beneficial.
- For simple datasets, a simple model will produce accurate results whilst also being computationally efficient, thus is preferred. Random Forest model goes against Ockams Razor theory and incorporates as much as possible into the system.
- On the other hand for a dataset with high dimensionality, implementing a Linear Regression model would require feature engineering and other data preprocessing steps such as data normalisation, data transformation, handling of missing values and outliers. In contrast these initial steps may not be necessarily required in a Random Forest model. It is therefore important to understand the dataset and select the most reasonable analytical approach. The lesson to take from this is that "there is no free lunch".
- To save time and computational resources, it is best to first implement a simple model on the most significant predicting variables. If the result is positive, it may not be necessary to fit a complex, expensive and hard-to-interpret model since the result may not be significantly improved.

Future work

- We acknowledge from the research work carried out that feeding additional observations into a Random Forest regression model improves performance. This is not the case with linear regression - i.e. the model does not benefit from the addition of more data.
- Generally with more training samples, there is more knowledge for RF to learn and so the model can subsequently capture the non-linearity of the structural data. Also, the model will improve with more appropriate features which can create a higher probability of choosing an optimum splitting feature [2].
- Taking the above into account we can say that in terms of future work the Random Forest method will be superior.

[1] Heysem Kaya, Pınar Tüfekci , Sadık Fikret Gürgen. Local and Global Learning Methods for Predicting Power of a Combined Gas & Steam Turbine, Proceedings of the International Conference on Emerging Trends in Computer and Electronics Engineering (2012) ICETCEE 2012, pp. 13-18

[2] Li, H., Leung, K.S., Wong, M.H. et al. BMC Bioinformatics (2014) 15: 291.

[3] Leo Breiman, 'Random Forests', Machine Learning 45, no. 1 (2001) 5-32.

[4] Paul Smith, Siva Ganesh, Ping Liu. A Comparison of Random Forest Regression and Multiple Linear Regression for Prediction in Neuroscience. Journal of neuroscience methods (2013). 220:10.1016

[5] Oshiro, T.M., Perez, P.S. and Baranauskas, J.A., 2012, July. How many trees in a random forest?. In MLDM (pp. 154-168).

[6] Alice Zheng. Evaluating Machine Learning Models. O'Reilly.

[7] James N. Morgan, John A. Sonquist. Problems in the Analysis of Survey Data, and a Proposal. Journal of the American Statistical Association (1963). Vol. 58 , Iss. 302.

[8] Ulrike Grömping. Variable Importance Assessment in Regression: Linear Regression versus Random Forest. The American Statistician 63, no. 4 (2009): 308-19.

[9] Ned Horning. Remote Sensing for Ecology and Conservation: A Handbook of Techniques.