

NAME:SATHEESH M

REG NO:422221104036

COLLEGE CODE :4222

COLLEGE NAME:SURYA GROUP OF INSTITUTIONS

Random forest

Random Forest is a supervised learning algorithm that employs ensemble learning method for classification and regression. It runs n number of regression trees, and combines them into a single model to make more accurate prediction than one single tree. RF constructs many decision trees at training, and predictions from all trees are combined to make the final prediction. Employing random sampling with replacement (bagging in machine learning terminology), RF helps data scientists reduce the variance associated with those algorithms that have high variance, typically decision trees. Given a training set of feature X and output Y , bagging repeatedly selects a random sample of the training set for β times ($b=1,2,\dots,\beta$) $\phi=1,2,\dots,\phi$ and fits the trees to these samples. For every tree, we obtain a sequence of instances which are randomly sampled replacement from the training set. Each sequence of instances corresponds to a random vector ϕ_k forming a specific tree. Since all the sequences will not be exactly the same, the decision trees constructed from them will also be slightly variant. De Aquino Afonso et al. suggest that the prediction of the K -th tree for an input X can be represented

:

$$h_k(X) = h(X, \phi_k), \forall k \in \{1, 2, \dots, K\} \quad h_{\phi} = h(X, \phi), \forall \phi \in \{1, 2, \dots, \phi\} \quad (9)$$

where K is the number of trees. As a tree splits, each of which randomly selects features to avoid correlations among features. Al paydin points out that a node S can be split into two subsets, S_1 and S_2 by selecting a threshold c that minimises the difference in the sum of squared errors..

$$SSE = \left(\sum_{i \in S_1} (v_i - 1/|S_1| \sum_{i \in S_1} v_i)^2 + \sum_{i \in S_2} (v_i - 1/|S_2| \sum_{i \in S_2} v_i)^2 \right) \quad \text{---} \sum_{i \in S_1} v_i^2 - 1/|S_1| \left(\sum_{i \in S_1} v_i \right)^2 + \sum_{i \in S_2} v_i^2 - 1/|S_2| \left(\sum_{i \in S_2} v_i \right)^2 \quad (10)$$

By following the same decision rules, we can predict any subtree as the mean or median output of instances. Finally, we can obtain the final prediction as an average of each tree's output.

:

$$h(X) = \frac{1}{K} \sum_{k=1}^K h_k(X) \quad h_{\phi} = \frac{1}{\phi} \sum_{\phi=1}^{\phi} h_{\phi}$$

EVALUATION

In machine learning, it is a general practice to rely on a correlation matrix to decide what features to be incorporated into our models. Table 3 presents the correlations between LRP_i and each feature, and shows that housing floor area has the largest correlation with prices, followed by property age, travelling time to Central District and floor level. In comparison, correlations between each orientation and prices are very close to zero and so they are excluded from our estimation. Then, we present the data visualisation to portray the relationship between property prices and each selected feature in figure 3. All the charts exhibit the expected relationship between the features and the housing prices.

Figure 3. Data visualisation

