

Challenge 4 Report

Kiyeob Lee, Amrita Veerabagu, Harinath

I. INTRODUCTION

In challenge4, we studied two topics: Linear Regression and Decision Tree. Both Linear Regression and Decision Tree belong to the supervised learning and a setup for the models is the following: Given a training dataset $D_n = \{(x_1, y_1), \dots, (x_n, y_n) \mid n \in \mathbb{N}\}$ where $x_i \in \mathbb{R}^d$ is an instance of the dataset and y_i is the corresponding true known label, our objective is to predict unknown \hat{y}_{n+1} provided a new data instance x_{n+1} comes in to the system. Linear regression is a model that starts with an assumption that there is a relationship between x_i and y_i . Decision Tree's objective is, similarly, given a training set, to learn which tree minimizes the misclassification error(or maximize the correct-classification). One way to implement decision tree is to use entropy method that you identify a classifier that maximizes an 'information gain' as well as to use GeNie method as learned in a tutorial.

II. LINEAR REGRESSION

One way to measure the misclassification is to use the squared-loss function, namely,

$$\mathcal{L}(h, (x, y)) = (h(x) - y)^2 \quad (1)$$

where $h(x)$ is a mapping of x into y and \mathcal{L} is a penalty of difference between $h(x)$ and y . From this loss function, we define the Mean Squared Error(MSE), namely,

$$\mathcal{L}(h) = \frac{1}{m} \sum_{i=1}^m (h(x_i) - y_i)^2 \quad (2)$$

Since we started with an assumption that there is a relationship between x_i and y_i and we have chosen the error is L2 norm difference, we can formulate the problem as a minimization problem:

$$\arg \min_w \mathcal{L}(h_w) = \arg \min_w \frac{1}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i)^2 \quad (3)$$

where $h(x_i)$ in equation (2) is replaced with an inner product by linearity assumption. To solve the problem, we can get the gradient of the objective function to be zero. That is,

$$\frac{2}{m} \sum_{i=1}^m (\langle w, x_i \rangle - y_i) x_i = 0 \quad (4)$$

which we can rewrite as $Aw = b$ where

$$A = \left(\sum_{i=1}^m x_i x_i^T \right) \quad \text{and} \quad b = \sum_{i=1}^m y_i x_i \quad (5)$$

If A is invertible, then the solution to the Empirical Risk Minimizer(ERM) problem is

$$w = A^{-1}b \quad (6)$$

If A is not invertible, it is still solvable, under a linearity assumption, because b is in the range of A .

III. DECISION TREE

The Classification and Regression Tree (CART) builds a multiway tree, finding for each node the categorical feature that will yield the largest information gain for categorical targets. This is an example for a greedy algorithm. The mathematical formulation of the binary tree is given below.

Given a training vector $x_i \in R_n, i = 1, \dots, l$ and its corresponding label value $y \in R_l$, a decision tree recursively partitions the space such that the samples are grouped together.

Let the data at node m be represented by Q . For each candidate split $\theta = (j, t_m)$ consisting of a feature j and threshold t_m , partition the data into $Q_{left}(\theta)$ and $Q_{right}(\theta)$ subsets,

$$\begin{aligned} Q_{left}(\theta) &= (x, y) | x_j \leq t_m \\ Q_{right}(\theta) &= Q / Q_{left}(\theta) \end{aligned}$$

The impurity at m is computed using an impurity function $H()$, the choice of which depends on the task being solved (classification or regression). Here, we use the $H()$ corresponding to classification, given by

$$H(X_m) = \sum_k p_{mk}(1 - p_{mk})$$

where p_{mk} is the proportion of class k observations in node m and is given by,

$$p_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k)$$

In the above equation, we consider that the target is a classification outcome taking on values $0, 1, \dots, K-1$, for node m , representing a region R_m with N_m observations.

Therefore the impurity at node m is given by,

$$G(Q, \theta) = \frac{n_{left}}{N_m} H(Q_{left}(\theta)) + \frac{n_{right}}{N_m} H(Q_{right}(\theta))$$

Thus we select the parameters which minimises the impurity,

$$\theta^* = \operatorname{argmin}_{\theta} G(Q, \theta)$$

Recurse for subsets $Q_{left}(\theta)$ and $Q_{right}(\theta)$ until the maximum allowable depth is reached, $N_m < \min_{samples}$ or $N_m = 1$.

IV. OBSERVATIONS PART 1-3

- 1) **Part 1** : In part 1, we built and predicted the linear regression model on same type of dataset i.e Trained model on White-Wine training set and predicted it on White-Wine testing with RMSD of 0.71 on test set. Similarly in case of Red-wine dataset, obtained a RMSD of 0.674 on test set. Model fit (RSq) of White wine - 0.984 and Red Wine - 0.987. So we can see that both model fit and test RMSD are good when model is trained and predicted on same type of dataset. Please look at the table for coefficient values of features from individual regression models.
- 2) **Part 2** : In part 2, decision tree model was built based on 11 input features and output is type of wine. Obtained an accuracy of 0.9953 on test set.
- 3) **Part 3** : In part 3, we first divided both red and white wine datasets in the ratio 80% and 20% as its training and test set. Then we trained two linear regression models using red-wine and white-wine training sets, Model1 and Model2 respectively. In order to test how well the models are performing on a different test set, the red-wine model was tested on white-wine test set and white-wine model was tested on red-wine dataset. The Model 1 gave an MSE of 0.401 and 0.729 in the red-wine and white-wine test set respectively. The Model 2 gave an MSE of 0.941 and 0.616 in the red-wine and white-wine test set respectively.

Feature	Relative Importance - Decision tree	Model Alpha - White Wine	Model Alpha - Red Wine
Total sulfur oxide	0.67	-0.0009	-0.0039
Chlorides	0.21	-1.157	-1.8033
Volatile acidity	0.06	-1.937	-1.0278
Sulphates	0.024	0.4167	0.9907
Density	0.017	1.9986	4.3731
pH	0.008	0.214	-0.5108
Alcohol	0.005	0.358	0.2931
Residual sugar	0.003	0.0236	0.0046
Fixed acidity	0.002	-0.052	0.0057
Free sulfur oxide	0.001	0.0046	0.0066
Citric acid	0.0001	-0.0189	-0.0732

V. OBSERVATIONS PART 4

- 1) From part 2 results , we can see that `DecisionTreeClassifier` works very well on both training and test set : which means given set of features can classify the data into wine type with great accuracy (0.9952). This means that data is well separable with respect to set of input features.
- 2) We can see the features being ranked based on importance values : The importance of a feature is computed as the (normalized) total reduction of the criterion brought by that feature. It is also known as the Gini importance. So in the table we see that Total Sulfur Oxide and Chlorides are top two features which help classify the data.
- 3) For these features we have also plotted the box plots with respect to wine type. (0- White, 1- Red). We can observe that box plots for 'chlorides' and 'total sulfur oxide' looks completely different for Red and White wine types. Now

this indicates why they were picked by Decision tree classifier as top features. Also we see that their ranges are also different among the types.

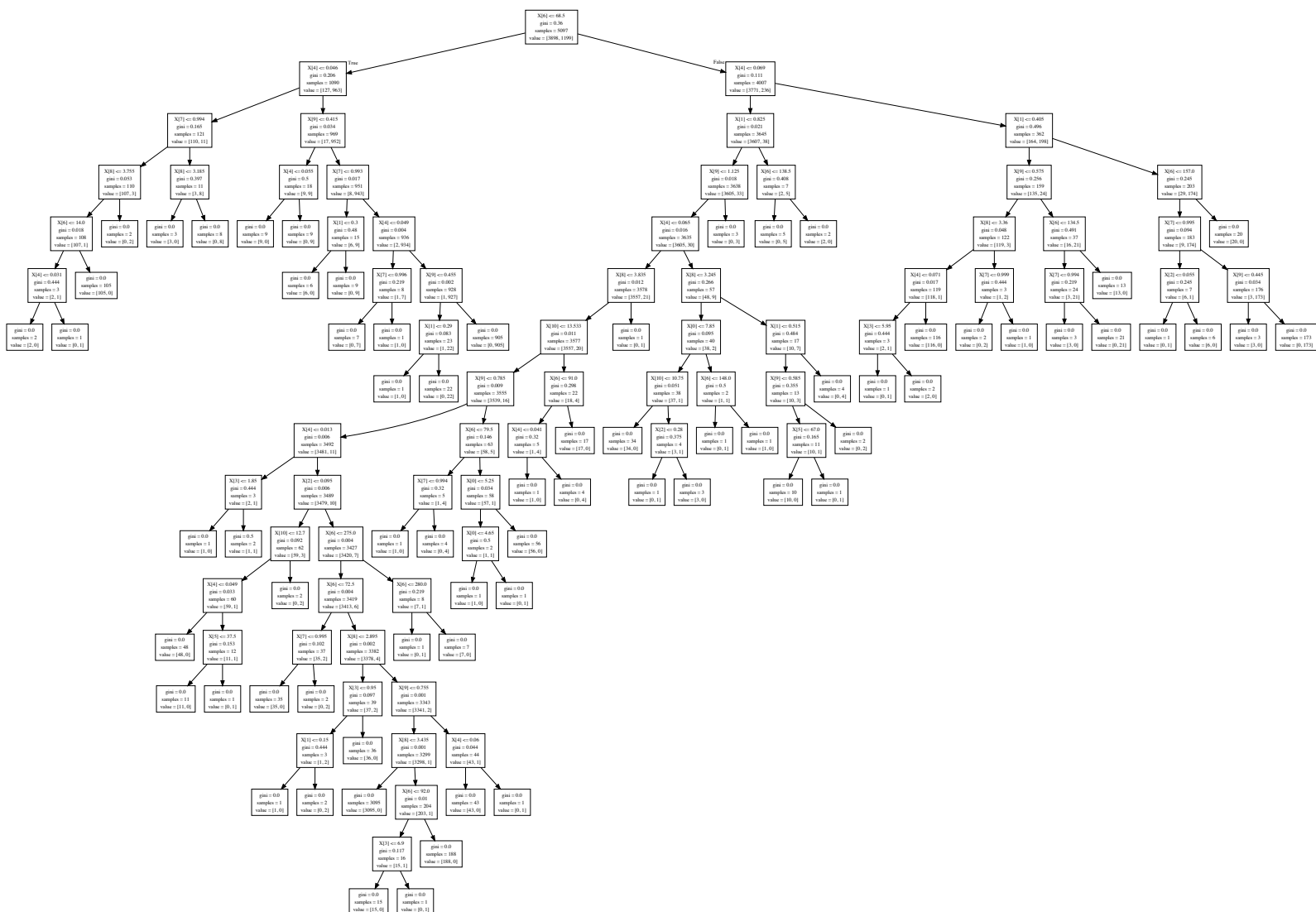
- 4) We can observe similar changes in range of features among types in case of 'pH', 'density', 'fixed acidity', 'volatile acidity' etc. This means the combination of features space in case of White Wine is different than feature space of Red Wine! And this is why we have good result in DecisionTreeClassifier.

5)

$$\hat{y}_i = \alpha_0 + \sum_j \alpha_i x_{i,j}. \quad (7)$$

Now in the case of regression ; as seen in the equation above each feature has their coefficient value $\{\alpha_i\}$ found out by Least squares estimate. We can see the alpha values for each features tabulated along with decision tree importance results.

- 6) Interesting observation is that : when we train model on White-Wine dataset and try to predict red-wine data we see that RMSD is bad compared to RMSD on Red-wine on Red-wine. Similarly we see that RMSD while predicting white dataset using model trained on Red is poor with respect to model trained on White dataset.
- 7) The reason being, when you train model on white wine dataset , compute alpha values and use them to predict red wine dataset (which has different ranges of X compared to white wine) predictions will be slightly off range since take product of α_i with X_i while predicting it.
- 8) From decision tree classifier results and box plots we can infer that input space of features in case of Red White is well separated. Hence when alpha values from one type of regression model is multiplied with completely different range of X values (extrapolation) usually doesn't guarantee good predictions.



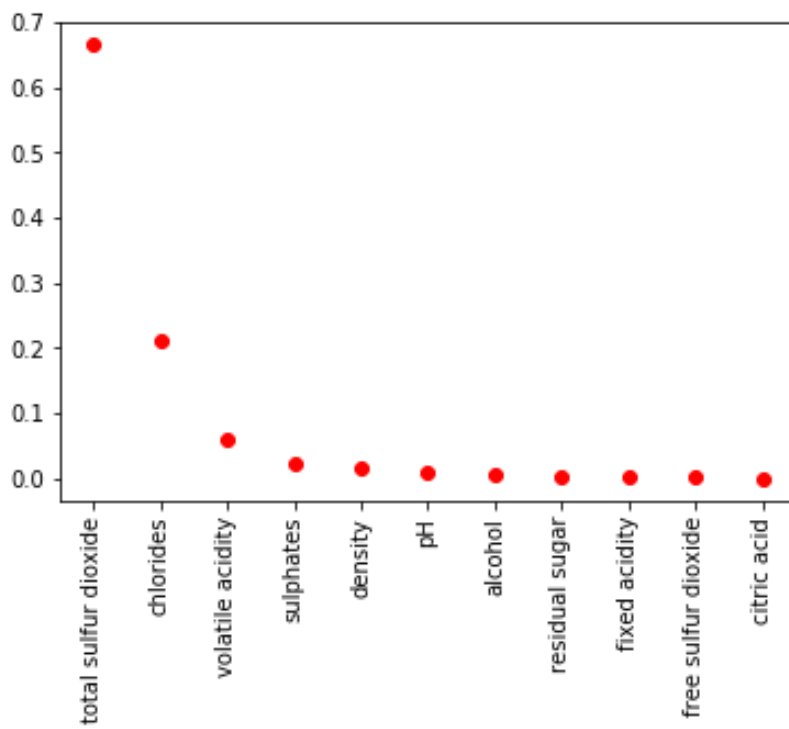


Fig. 1: Feature importance values of decision tree classifier

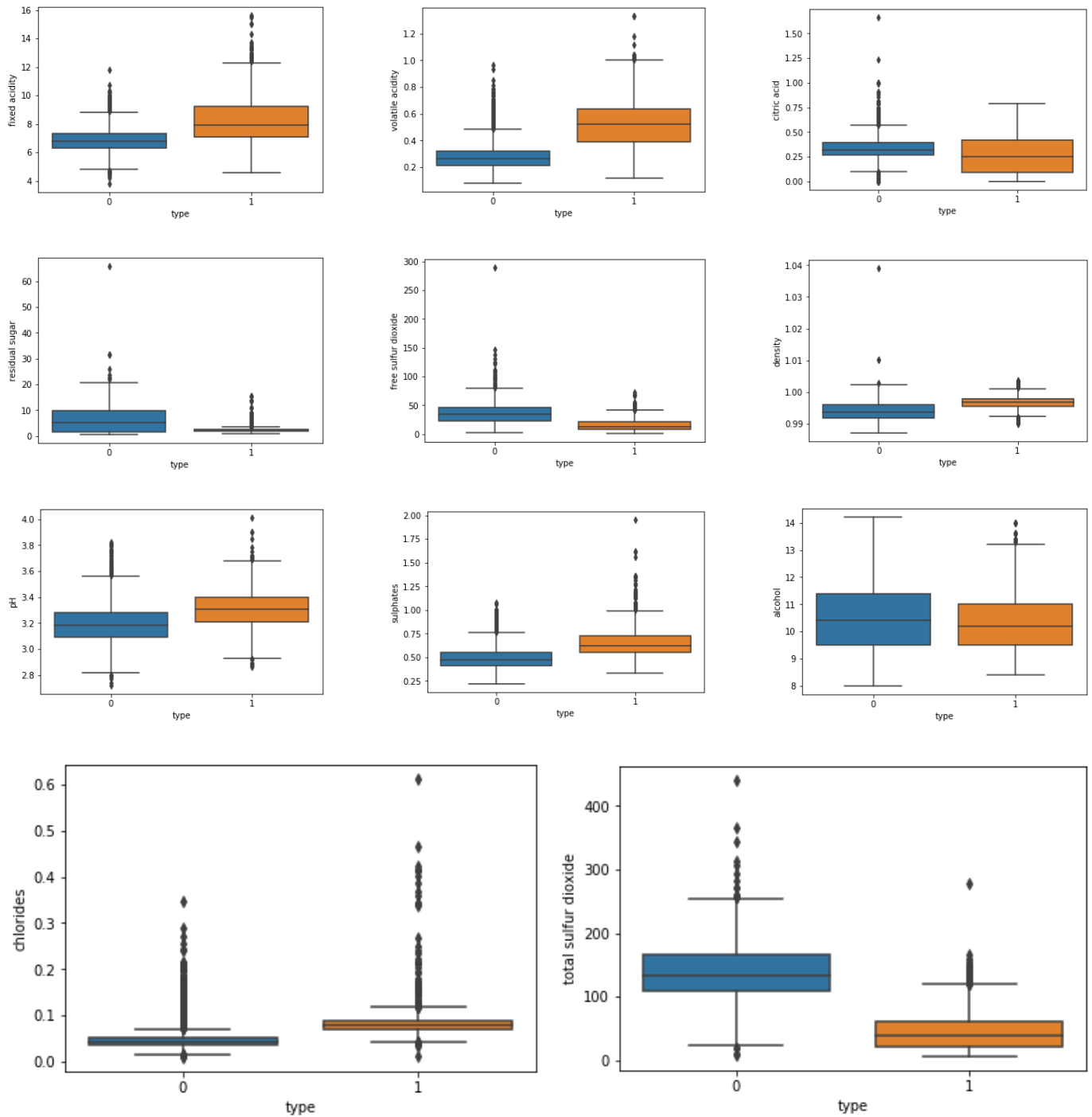


Fig. 2: Box Plots of features for each wine type