

1. Decision Trees

(a) overall entropy: $-\sum_{i=1}^C p(x_i) \log_b q(x_i)$

$$S = 5 \quad S' = 5$$

$$-\gamma_2 \log_2 (\gamma_2) = \gamma_2$$

$$-\sum_{i=1}^2 p(x_i) \log_2 = \gamma_2 + \gamma_2 = 1$$

$$IB(T, A) = \text{Entropy}(T) - \sum_{j \in A} \frac{|T_j|}{|T|} \cdot \text{Entropy}(T_j)$$

a Color:

$$\begin{array}{c} R \\ S | S' \\ \hline 3 | 2 \end{array}$$

$$\text{entropy} = -\left(\frac{3}{5} \log_2 \left(\frac{3}{5}\right) + \frac{2}{5} \log_2 \left(\frac{2}{5}\right)\right) = 0.97$$

$$\begin{array}{c} Y \\ S | S' \\ \hline 2 | 3 \end{array}$$

$$\text{entropy} = -\left(\frac{2}{5} \log_2 \left(\frac{2}{5}\right) + \frac{3}{5} \log_2 \left(\frac{3}{5}\right)\right) = 0.97$$

$$IB = 1 - (.5(0.97) + 0.5(0.97)) \approx \underline{0.03}$$

• T-type

sports $\frac{5}{4} \left| \begin{smallmatrix} 5 \\ 2 \end{smallmatrix} \right.$ $e_1 = -\left(\frac{4}{6} \log_2\left(\frac{4}{6}\right) + \frac{2}{6} \log_2\left(\frac{2}{6}\right)\right) \approx 0.9183$

SUV $\frac{5}{1} \left| \begin{smallmatrix} 5 \\ 3 \end{smallmatrix} \right.$ $e_2 = -\left(\frac{4}{4} \log_2\left(\frac{4}{4}\right) + \frac{3}{4} \log_2\left(\frac{3}{4}\right)\right) \approx 0.8113$

$$\frac{6}{10}e_1 + \frac{4}{10}e_2 \approx 0.55098$$

$$I_6 \approx 1 - 0.55098 \approx \underline{0.44902}$$

• Origin

Domestic

$$\frac{5}{2} \left| \begin{smallmatrix} 5 \\ 3 \end{smallmatrix} \right. e_1 = e_2 = -\left(\frac{2}{5} \log_2\left(\frac{2}{5}\right) + \frac{3}{5} \log_2\left(\frac{3}{5}\right)\right)$$

imported

$$\frac{5}{3} \left| \begin{smallmatrix} 5 \\ 2 \end{smallmatrix} \right. \text{(same as Color split)}$$

$$I_6 \approx \underline{0.03}$$

Best feature for first split is "T-type" as it gave the most I₆

(b) split on type for the first branch.
For each subsequent branch, repeat the
same process as in (a) for each
feature not yet used in that
branch. Then calculate
the most likely feature with
information gain or entropy
until branches are
reached.

(c) We can stop when the number of
nodes in a node is less than some
predetermined limit, the purity of the node
is more than some predetermined
(limit), we have already split on the
three features in every branch to
a leaf node or stop when the predicted
values for all records are identical.

(d) May have \exists no need in most cases to standardize the features when using a decision tree. Monotonic transforms do not affect the decision tree splits.

(e) Decision trees are robust to outliers as they split data into multiple unique groups rather than drawing a hyperplane or single weighted boundary.

2. True or False, Simple Explanations

- (a) False, Boosting uses weak learners with the idea of producing many cheaper weak learners instead of a few expensive strong learners.
- (b) False, in Bagging, the algorithm can look through all features to find optimal split point, while random forest is limited to random sample of features of which to search.
- (c) While both bagging and boosting decrease variance and tend to a model with high stability, only boosting can reduce bias. Also, bagging can be efficient in reducing overfitting, but boosting can increase overfitting and requires careful tuning of the hyper-parameters. Boosting can however support different loss functions.

(d) with certain activation functions such as the sigmoid function, the input is squared, so changes in the input lead to a small change in output. With many layers, this can cause a vanishing gradient in back propagation. The ReLU is discontinuous, and the gradient is either $\underline{0}$ or $\underline{1}$, so gradients can't get squared. However, you can set ReLUs which return 0 forever, so never update them.

3. Overfitting Mitigation Strategies

- 3.1) Yes, a larger dataset most often is a better representation of population data, and less likely to generalize, so it will help avoid overfitting.
- 3.2) No, allowing the model to train for more iterations decreases bias and increases variance, as the model will be more fitted to the data and less generalizable with more iterations.
- 3.3) No, as more parameters leads to a more complex higher dimensional model, which is more likely to overfit.
- 3.4) Yes, using this dropout method, we make the model learn some but not all information in dropped nodes making it more generalizable and less prone to overfitting.

3.5) No, a specialized chip might be faster, but the process is the same and so is the risk of overfitting.

3.6) Yes, changing the initial values/strategy can lead to finding a local minimum instead of a global one, less overfitting as the model is less exact.

4) Principal Component Analysis

(a) The advantages of dimensionality reduction are a result of correlated features which lowers the input dimensions and improves algorithm performance. It also reduces overfitting and improves visualization methods.

However, independent variables become less readable and interpretable, data standardization must be done. Nevertheless, other dimensionality reduction will lead to terrible results, and dimensionality reduction will lead to information loss.

(b)

- Data with an inherently linear structure will work well with PCA, as we can reduce dimensionality by fully ge principal component along which the data lies and due to low variation of other factors, we can eliminate them.
- Data lying on a hyperbolic plane will not work well with PCA, as unlike the linear dataset, it is not always easy to center points in hyperbolic space. As such, the steps in PCA do not work well with hyperbolic data.
- Data containing unscaled features will not work with PCA, as without normality we cannot be certain which features have more impact than the others.
- Statistically independent features indicate each feature provides no information and should not be discarded, so PCA should not be used.

5) Neural Networks

(a) All activation functions share a similar property: non-linearity and ensure gradients remain large enough (do not vanish). Common choices are sigmoid, tanh, ReLU and softmax variants, softplus and swish. Input layers do not have activation functions, as the activation function is used in the forward propagation of data and the input is the first step. Hidden layers can use ReLU (\rightarrow), sigmoid (\curvearrowleft), and tanh (\downarrow). Sigmoid and tanh can lead to vanishing gradients, while ReLU prevents this. As such, it works well for the hidden nodes. The output activation function depends on the type of prediction that one is doing. If the target is binary, it would use a sigmoid function to get probability of class membership; if the β multiclass classification, then the output layer will have one node per class so a softmax activation should be used to help aggregate them. If it is a regression problem, just use a linear activation.

$$(b) \quad x_0 = -2 \quad x_1 = 1 \quad w_0 = 0.5 \quad w_1 = 0.75 \quad w_2 = -0.25 \\ \text{bias term} = 1$$

$$p = \text{ReLU}(w_0 x_1 + w_1 x_2 + w_2(\text{bias})) = \text{ReLU}(-1 + \frac{3}{4} + 1) = \text{ReLU}(\frac{3}{4}) = \frac{3}{4}$$

$$g = \text{ReLU}$$

$$\text{hidden 0} = \text{ReLU}(w_0 x_0 + w_1 x_1 + b) = \max(0, [-1 + \frac{3}{4} + 1]) = \frac{3}{4}$$

$$\text{hidden 1} = \text{ReLU}(w_0 x_0 + w_1 x_1 + b) = \frac{3}{4}$$

$$\text{hidden 2} = \sigma(w_0 x_0 + w_1 x_1 + b) = \frac{1}{1 + e^{-\frac{3}{4}}} \approx 0.6792$$

$$\text{output} = \frac{\frac{3}{4} + \frac{3}{4} - \gamma_1(0.6792)}{2} = \underline{0.5802}$$