

Fundamentals of Artificial Intelligence

Homework - 5
CSCE - 5210

Ajay Reddy kudumula
11580520

- Q1. Suppose, We have a System that observes a person's TV Watching habits in Order to recommend other TV Shows the person's may like. Suppose that we have a characterized each show by Whether it's a Comedy, Whether it's features doctors, Whether it features lawyers and Whether it has guns, Suppose We are given the examples of fig 7.23 about Whether person likes Various TV Shows.

Example	Comedy	Doctors	lawyers	Guns	likes
e ₁	false	true	false	false	false
e ₂	true	false	true	false	true
e ₃	false	false	true	true	true
e ₄	false	false	true	false	false
e ₅	false	false	false	true	false
e ₆	true	false	false	true	false
e ₇	true	false	false	false	true
e ₈	false	true	true	true	true
e ₉	false	true	true	false	false
e ₁₀	true	true	true	false	true
e ₁₁	true	true	false	true	false
e ₁₂	false	false	false	false	false

We want to use this dataset to learn the value of Likes (i.e., to predict which TV shows the person would like based on the attribute of TV show)

b) Do the same as in part (a) but with Sum of Squares error

Ans:- The optimal decision tree with one node predicts likes = false (or not like). It has 5 errors. Then the sum of square of error will be

$$5 \times \left(\frac{7}{12}\right)^2 + 7 \times \left(\frac{5}{12}\right)^2 = 1.70 + 1.915 = 2.92.$$

d) Do the same as in part (c), but with the sum of square of errors.

Ans:- An optimal solution with depth 2 is if $\text{likes} = \text{true}$ else $\text{like} = \text{false}$.

It has 3 errors. At root are all of examples (e_1, \dots, e_{12}) . Filtered to $\text{likes} = \text{true}$ node are $\{e_2, e_3, e_4, e_8, e_9, e_{10}\}$. Filtered to $\text{likes} = \text{false}$ node are $\{e_1, e_5, e_6, e_7, e_{11}, e_{12}\}$. If $\text{likes} = \text{true}$ then $\text{likes} = 4/6$ else $\text{like} = 3/6$.

The sum of square error is

$$2 \left(\frac{4}{6}\right)^2 + 4 \left(\frac{2}{6}\right)^2 + 7 \left(\frac{5}{6}\right)^2 + 5 \left(\frac{1}{6}\right)^2 = 2.16$$

e) What is the smallest tree that correctly classifies all training examples? Does a top-down decision tree that optimizes the information gain at each step represent the same function.

Ans The smallest decision tree is if guns then (if lawyers then likes = true else like = false)
else

(if Comedy then likes = true else likes = false)

One way to find such trees is to do a two-step look ahead for each property, check a Split on that property and then do a Split on each leaf before evaluating the Split.

f) Give two instance not appearing in the examples of fig F23 and show how they are classified using the Smallest decision trees use this to explain the bias inherit in the tree (How does the bias give you these particular predictions)

Ans: The logistic regression learns algorithm can learn any linearly Separately classification. The error can be made arbitrarily small for arbitrary sets of examples if and only if the target classification is linearly Separable.

Q.2) It is possible to define a regularizer to minimize $\sum_e (\text{error}_n(e) + \lambda * \text{regularizer}_n)$ rather than Formula 7.5. How this is different than the existing regularizer? [Hint: Think about this effects multiple datasets or for cross validation].

Suppose λ is a set by k-fold cross validation, and then the model is learned for the whole dataset. How would the algorithm be different from the original ways of defining a regularizer and this alternative way? [Hint: There is a different number of examples used for regularization than there is the full dataset; does this matter?] Which works better in practice?

Ans (I) The regularizer in formula 7.5 is designed to minimize the sum of the error for each data point plus a penalty term that encourages the model to be simpler (i.e., to have fewer parameters). This regularizer is effective at preventing overfitting on a single data set.

(II) The alternative regularizer is more effective at preventing overfitting on preventing multiple dataset or when using cross-validation for problems. Where we want the model to have more parameters may be a good choice.

Explanation:

(I) If you are working with multiple datasets or using cross-validation. This regularizer may not be ideal. In particular, if the model is fit on multiple datasets, the regularizer will encourage the model to be simpler. This may not be desirable.

if you want the model to be able to learn different patterns in each dataset.

⇒ Alternatively, you could define a regularizer that minimizes the sum of the errors. For each data point plus a penalty term that encourages the model to be more flexible. This regularities would be more effective at preventing over fitting on multiple datasets or when using cross-validation.

⇒ There are a few different way to define such a regularities one option is to use the ℓ_1 norm which encourages the model to have small parameter values. These are just two of many possible options.

⇒ In general, the choice of regularities will depend on the specific problem and data. There is no single best regularized for all problems. However, for more problems where you want the model to be able to learn different patterns in multiple datasets. A regularities that encourages the model to have more parameters may be a good choice.

(II) ⇒ There are few key differences between the original regularizes and the alternative regularities. First, the original regularizer is defined using a single dataset, while the alternative regularizer is defined using multiple datasets. This means that the alternative regularizer will be more effective at preventing over fitting on multiple datasets.

⇒ Second, the Original regularizer encourages the model to be simpler (i.e., to have fewer parameters). While the alternative regularizer encourages the model to be more flexible. This may be desirable if we want the model to be able to learn different patterns in each dataset.

⇒ Third, the Original regularizer is defined using the l_2 norm, while the Alternative regularizer is defined using l_1 norm. This means that the alternative regularizers will encourage the model to have fewer non-zero parameters.

⇒ Fourth, the Original regularizer is fit to the entire dataset, while, the Alternative regularizer is fit on the subset of dataset. This means that the alternative.

⇒ Regularities will be less likely to overfit on the data. Overall, the alternative regularizer is more effective at preventing overfitting on multiple datasets or when using cross-validation.

⇒ There is no single best regularities for all problems. The choice of regularities will depend on the specific problem and data in general. For problems where they want the model to be able to learn different patterns in multiple datasets, a regularizer that encourages the model to have more parameters may be a good choice.