
Data Innovators - Report for Assignment 2 - CS771

1. Kshitij Bhardwaj - 230580

2. Harsh Gupta - 230445

3. Priyanka Arora - 230799

4. Parnika Mittal - 230736

5. Harshit Agarwal - 230458

Abstract

This report describes the development of a decision tree model for word prediction using bigrams (adjacent character pairs) as features. The model aims to predict a single word based on a list of bigrams. The report outlines the design decisions taken for the model, including feature engineering, the decision tree algorithm, and hyperparameters.

1 Decision Tree Algorithm

1.1 Splitting Criterion - Entropy

1. Precision on Gini Criterion: 0.3783626862783046
2. Precision on Entropy Criterion: 0.4337139539384556

Why Entropy Might Be More Suitable for This Specific Bigram Prediction Problem?

Entropy: Measures the overall uncertainty or randomness in a dataset. Higher entropy indicates more uncertainty about the correct class (word) for a given bigram feature.

Gini Impurity: Focuses on the probability of misclassifying a sample if randomly drawn from a node. Lower Gini impurity implies a better separation between classes.

Reasoning for Better Performance with Entropy:

1. Bigram Ambiguity: Bigrams might be ambiguous and not uniquely identify a specific word. For example, the bigram "th" could appear in words like "the," "this," "thing," etc.
2. Entropy Captures Ambiguity: Entropy inherently penalizes uncertainty and aims to distribute samples across all possible classes (words) for a given bigram. This can be beneficial when dealing with ambiguous features like bigrams.

Imagine a dataset with bigrams "th" and "yi" and three words: "the" (2 occurrences), "thing" (1 occurrence), and "playing" (1 occurrence).

$$Gini = 1 - \sum_{i=1}^C (p_i)^2 \quad (1)$$

Gini Impurity Calculation:

- Node with bigram "th": Classes (words) - "the" (2) and "thing" (1).
- Gini impurity = $1 - (0.64^2 + 0.36^2) = 0.384$
- Node with bigram "yi": Class (word) - "playing" (1).
- Gini impurity = $1 - 1^2 = 0$

$$E(S) = \sum_{i=1}^c -p_i \log_2(p_i) \quad (2)$$

Entropy Calculation:

- Node with bigram "th": Entropy = $-(2/4 * \log_2(2/4) + 1/4 * \log_2(1/4)) = 0.811$
- Node with bigram "yi": Entropy = $-(1/1 * \log_2(1/1)) = 0$
- Overall Entropy (weighted average) might be slightly higher than Gini impurity.

1.2 Stopping Criterion

There are two main criteria that determine when to stop expanding a decision tree and make a node a leaf:

1. Purity of the Node : 'max_depth = None'

This criterion checks if all samples reaching a particular node in the tree belong to the same class (category). If so, the node is considered pure, and there's no need for further splitting. This is because the model can already make a confident prediction (the class) for any sample that reaches that node. To maximise this parameter usefulness, we set the 'max_depth = None'

If None, then nodes are expanded until all leaves are pure or until all leaves contain less than 'min_samples_split' samples.

2. Minimum Samples Per Leaf : 'min_samples_leaf = 1'

This hyperparameter sets a threshold for the minimum number of samples allowed in a leaf node. Expanding the tree further might be pointless if there are very few samples remaining in a node. With too few samples, the decision tree might learn overly specific rules that don't generalize well to unseen data (underfitting).

Parameter Value	1	2	4
Precision	0.4256	0.1661	0.0689

The Relative gain in Model Size Reduction and Time Taken, does not override the loss in precision.

1.3 Pruning : None

Our Problem statement does not require external Pruning strategies, it well fits the data provided to output a respectable accuracy.

1.4 Other Hyperparameters

1. **min_samples_split : 2**

With a higher min_samples_split, the overall dataset might have more unique bigrams and next words, potentially leading to a higher initial entropy (H) before splitting.

Although there are chances that Information Gain is better overall after 3 layers or so, the precision always decreases on increasing the value.

Parameter Value	2	3	4
Precision	0.4337	0.2721	0.2106

2 References

1. scikit-learn DecisionTreeClassifier Documentation
2. Entropy, Info Gain, Gini Index and CCP Pruning Article
3. Article on Classification Trees Pruning Optimization with GridSearchCV and Cost Complexity Fn.