Learning with Decision Trees

CS771: Introduction to Machine Learning Piyush Rai

Announcement

- An extra office hour every week (Saturday 3-4pm)
 - Different from my regular in-person office hours (Wed, 6-7pm)
- To be held online (Google Meet: <u>https://meet.google.com/dup-mozx-swe</u>)
- Can ask doubts etc from the previous classes



Decision Trees

• A Decision Tree (DT) defines a hierarchy of rules to make a prediction



- Root and internal nodes test rules. Leaf nodes make predictions
- DT learning is about learning such a tree from labeled training data



Decision Tree Learning: The Basic Idea

Recursively partition training data till you get (roughly) homogeneous regions



- Some typical prediction rules for each region
 - Use a constant label (e.g., majority) if region fully/almost homogeneous
 - Learn another prediction model (e.g., LwP) if region not fully homogeneous



Decision Tree for Classification: An Example







DT is very efficient at test time: To predict the label of a test point, nearest neighbors will require computing distances from 48 training inputs. DT predicts the label by doing just 2 feature-value comparisons! Way more fast!!! Remember: Root node contains all training inputs. Internal/leaf nodes receive a subset of training inputs



Decision Tree for Regression: An Example





To predict the output for a test point, nearest neighbors will require computing distances from 15 training inputs. DT predicts the label by doing just at most 2 feature-value comparisons! Way more fast!!!



How to decide which rules to YES NO $x_1 > 3.5$ YES NO NO YES $x_2 > 3?$ $x_2 > 2$? Predict Predict Predict Predict Red Red Green Green In general, constructing DT is an 5 6 The rules are organized in the Feature 1 (x_1) DT such that most informative Hmm.. So DTs are like rules are tested first the "20 questions" game (ask the most Informativeness of a rule is of related useful questions first) to the extent of the purity of the split arising due to that rule. More informative rules yield more pure splits

Constructing a Decision Tree

Feature 2 (x_2) **C**

Given some training data, what's the "optimal" DT?



test for and in what order?

How to assess informativeness of a rule?

intractable problem (NP-hard)



CS771: Intro to ML

Often we can use some "greedy" heuristics to construct a "good" DT

To do so, we use the training data to figure out which rules should be tested at each node

The same rules will be applied on the test inputs to route them along the tree until they reach some leaf node where the prediction is made



Techniques to Split at Internal Nodes?

- This decision/split can be done using various ways, e.g.,
 - Testing the value of a single feature at a time (such internal node called "Decision Stump")

With this approach, all features (2 real-valued features in this example) and all possible values of each feature need to be evaluated in selecting the feature to be tested at each internal node. If features binary/discrete (only finite possible values), it is reasonably easy





DT methods based on testing a single feature at each internal node are faster and more popular (e.g., ID3, C4.5 algos)



- Testing the value of a combination of features (maybe 2-3 features)
- Learning a classifier (e.g., LwP or some more sophisticated classifier)



DT methods based on learning and using a separate classifier at each internal node are less common. But this approach can be very powerful and sometimes used in some advanced DT methods



Internal Nodes: Good vs Bad Splits

- Recall that each internal node receives a subset of all the training inputs
- Regardless of the criterion, the split should result in as "pure" groups as possible
 - Meaning: After split, in each group, majority of the inputs have the same label/output



- For classification problems (discrete outputs), entropy is a measure of purity
 - Low entropy \Rightarrow high purity (less uniform label distribution)
 - Splits that give the largest reduction (before split vs after split) in entropy are preferred (this reduction is also known as "information gain")

Entropy and Information Gain

- Assume a set of labelled inputs S from C classes, p_c as fraction of class c inputs
- Entropy of the set **S** is defined as $H(S) = -\sum_{c \in C} p_c \log p_c$
- Suppose a rule splits $m{S}$ into two smaller disjoint sets $m{S}_1$ and $m{S}_2$
- Reduction in entropy after the split is called information gain



CS771: Intro to ML

11

Uniform sets (all classes

roughly equally present)

sets low

have high entropy; skewed

Decision Tree for Classification: Another Example

- Deciding whether to play or not to play Tennis on a Saturday
- Each input (Saturday) has 4 categorical features: Outlook, Temp., Humidity, Wind
- A binary classification problem (play vs no-play)
- Below Left: Training data, Below Right: A decision tree constructed using this data

day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no



Entropy and Information Gain

- Let's use IG based criterion to construct a DT for the Tennis example
- At root node, let's compute IG of each of the 4 features
- Consider feature "wind". Root contains <u>all</u> examples S = [9+,5-]

$$H(S) = -(9/14) \log_{2}(9/14) - (5/14) \log_{2}(5/14) = 0.94$$

$$S_{\text{weak}} = [6+, 2-] \Rightarrow H(S_{\text{weak}}) = 0.811$$

$$S_{\text{strong}} = [3+, 3-] \Rightarrow H(S_{\text{strong}}) = 1$$

$$IG(S, wind) = H(S) - \frac{|S_{\text{weak}}|}{|S|} H(S_{\text{weak}}) - \frac{|S_{\text{strong}}|}{|S|} H(S_{\text{strong}}) = 0.94 - 8/14 * 0.811 - 6/14 * 1 = 0.048$$

dav

2

3

4

outlook

sunny

sunny

overcast

rain

rain

temperature

hot

hot

hot

mild

cool

humidity

high

high

high

high

normal

normal

wind

weak

strong

weak

weak

weak

strong

play

no

no

yes

yes

- Likewise, at root: IG(S, outlook) = 0.246, IG(S, humidity) = 0.151, IG(S, temp) = 0.029
- Thus we choose "outlook" feature to be tested at the root node
- Now how to grow the DT, i.e., what to do at the next level? Which feature to test next?
- Rule: Iterate for each child node, select the feature with the highest IG

		TLUT		
\$771:	In	tro	to	MI

Growing the tree

day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no



- Proceeding as before, for level 2, left node, we can verify that
 - IG(S,temp) = 0.570, IG(S, humidity) = 0.970, IG(S, wind) = 0.019
- Thus humidity chosen as the feature to be tested at level 2, left node
- No need to expand the middle node (already "pure" all "yes" training examples ③)
- Can also verify that wind has the largest IG for the right node
- Note: If a feature has already been tested along a path earlier, we don't consider it again

14

When to stop growing the tree?



Stop expanding a node further (i.e., make it a leaf node) when

- It consist of all training examples having the same label (the node becomes "pure")
- We run out of features to test along the path to that node
- The DT starts to overfit (can be checked by monitoring the validation set accuracy)
 To help prevent the tree from growing too much!

Important: No need to obsess with too much for purity

- It is okay to have a leaf node that is not fully pure, e.g., this
- At test inputs that reach an impure leaf, can predict probability of belonging to each class (in above example, p(red) = 3/8, p(green) = 5/8), or simply predict the majority label



Wind

Weak

Avoiding Overfitting in DTs

- Desired: a DT that is not too big in size, yet fits the training data reasonably
- Note: An example of a very simple DT is "decision-stump"
 - A decision-stump only tests the value of a single feature (or a simple rule)
 - Not very powerful in itself but often used in large ensembles of decision stumps
- Mainly two approaches to prune a complex DT
 - Prune while building the tree (stopping early) Either can be don
 - Prune after building the tree (post-pruning)
- Criteria for judging which nodes could potentially be pruned
 - Use a validation set (separate from the training set)
 - Prune each possible node that doesn't hurt the accuracy on the validation set
 - Greedily remove the node that improves the validation accuracy the most
 - Stop when the validation set accuracy starts worsening
 - Use model complexity control, such as Minimum Description Length (will see later)





Decision Trees: Some Comments

- Gini-index defined as $\sum_{c=1}^{C} p_c (1 p_c)$ can be an alternative to IG
- For DT regression¹, variance in the outputs can be used to assess purity

For regression, outputs are real-valued and we don't have a "set" of classes, so quantities like entropy/IG/gini etc. are undefined

- When features are real-valued (no finite possible values to try), things are a bit more tricky
 - Can use tests based on thresholding feature values (recall our synthetic data examples)
 - Need to be careful w.r.t. number of threshold points, how fine each range is, etc.
- More sophisticated decision rules at the internal nodes can also be used
 - Basically, need some rule that splits inputs at an internal node into homogeneous groups
 - The rule can even be a machine learning classification algo (e.g., LwP or a deep learner)
 - However, in DTs, we want the tests to be fast so single feature based rules are preferred
- Need to take care handling training or test inputs that have some features missing

Ensemble of Trees

- Ensemble is a collection of models
- Each model makes a prediction. Take their majority as the final prediction Each tree is trained on a
- Ensemble of trees is a collect of simple DTs
 - Often preferred as compared to a single massive, complicated tree <u>inputs/features</u>
- A popular example: Random Forest (RF)

An RF with 3 simple trees. The majority prediction will be the final prediction



subset of the training

All trees can be

trained in parallel



- XGBoost is another popular ensemble of trees
 - Based on the idea of "boosting" (will study boosting later) simple trees
 - Sequentially trains a set of trees with each correcting errors of previous ones



Decision Trees: A Summary

Some key strengths:

- Simple and easy to interpret
- Nice example of "divide and conquer" paradigm in machine learning
- Easily handle different types of features (real, categorical, etc.)
- Very fast at test time
- Multiple simple DTs can be combined via ensemble methods: more powerful
- Used in several real-world ML applications, e.g., recommender systems, gaming (Kinect)

Some key weaknesses:

Learning optimal DT is (NP-hard) intractable. Existing algos mostly greedy heuristics

re 2 (x₂)

eatur V

Can sometimes become very complex unless some pruning is applied



5

Feature 1 (x_1)

6

