Applied Machine Learning

Decision Trees

Reihaneh Rabbany



COMP 551 (winter 2023) 1

Motivation

What we have left to cover for this course:

Classification and regression trees Linear support vector machines Bagging & boosting Unsupervised learning Dimensionality reduction

> from 2020 Kaggle's survey on the state of Machine Learning and Data Science, you can read the full version here



METHODS AND ALGORITHMS USAGE

Learning objectives

Decision trees:

- how does it model the data?
- how to specify the best model using a cost function
- how the cost function is optimized

Decision trees: motivation



pros.

- decision trees are interpretable
- they are not very sensitive to outliers
- do not need data normalization

cons.

• they could easily overfit and are unstable to small changes in input data

image credit:https://mymodernmet.com/the-30second-rule-a-decision/

					-
	<tumorsiz< th=""><th>ze, texture, p</th><th>erimeter> ,</th><th><cancer></cancer></th><th></th></tumorsiz<>	ze, texture, p	erimeter> ,	<cancer></cancer>	
$x^{(1)}$	<18.2,	27.6,	117.5> ,	< No >	$y^{(1)}$
$x^{(2)}$	<17.9,	10.3,	122.8> ,	< No >	$y^{(2)}$
$x^{(3)}$	<20.2,	14.3,	111.2> ,	< Yes >	$y^{(3)}$
			•		
•			•		:
$x^{(N)}$	<15.5,	15.2,	135.5> ,	< No >	$y^{(N)}$

Notation overview

Our datasets consists of N pairs of input vector and corresponding target or label $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$

we use N to denote the size of the dataset and $\frac{n}{n}$ for indexing a pair of input, label x, y denote an input and label pair where

x is a *D*-dimensional vector: $x = [x_1, x_2, \dots, x_D]$

we use *D* to denote the number of *features* (dimensionality of the input space)

 $y \in \{1,\ldots,C\}$ for classification problems, we use ${\scriptscriptstyle C}$ for number of classes

 $x_d^{(n)}$ $n \in [1 \dots N]$ indexes an instance, row index, e.g. which patient $d \in [1 \dots D]$ indexes a feature, column index, e.g. which measurement

e.g.
$$x_2^{(3)}$$

Decision trees: idea

- divide the input space into regions $\mathbb{R}_1, \ldots, \mathbb{R}_K$ using a tree structure
- assign a prediction **label** to each region

for classification this is the class label for regression, this is a real scalar or vector

$$f(x) = \sum_k w_k \mathbb{I}(x \in \mathbb{R}_k)$$

how to build the regions and the tree?

split regions successively based on the value of a single variable called test

each region is a set of conditions $\mathbb{R}_2 = \{x_1 \ge t_1, x_2 \le t_4\}$





Prediction per region

What constant w_k should we use for prediction in each region \mathbb{R}_k ?

Classification

count the frequency of classes per region, predict the most frequent label or return $w_k = ext{mode}(y^{(n)} | x^{(n)} \in \mathbb{R}_k)$ probability

Regression

use the mean value of training data-points in that region $w_k = ext{mean}(y^{(n)}|x^{(n)} \in \mathbb{R}_k)$

example: predicting survival in titanic

is sex male?

no





Possible tests

next questions: what are all the possible tests? which test do we choose next?



Continuous features

all the values that appear in the dataset can be used to split

Categorical features

if a feature can take C values $\ x_i \in \{1,\ldots,C\}$

convert that feature into C binary features (one-hot coding) $x_{i,1},\ldots,x_{i,C}\in\{0,1\}$

split based on the value of a binary feature

alternatives:

- **multi-way** split: can lead to regions with few datapoints
- binary splits that produce balanced subsets



Food Name		Categorical #		Calories	
Apple		1		95	
Chicken		2		231	
Broccoli		3		50	
Apple	Chi	cken	Broccoli		Calories
1	0		0	0	
0	1		0		231
0	0		1		50

Cost function

find a **decision tree** minimizing the following cost function, this cost function specifies "what is a good decision or regression tree?"

first calculate cost per region \mathbb{R}_k

Similar to other ML algorithms minimize a cost function or maximize an objective function



regression cost

we predict $w_k = ext{mean}(y^{(n)} | x^{(n)} \in \mathbb{R}_k)$ for region \mathbb{R}_k

mean squared error (MSE)

 $\mathrm{cost}(\mathcal{D}) = \sum_k rac{N_k}{N} \mathrm{cost}(\mathbb{R}_k, \mathcal{D})$

$$\mathrm{cost}(\mathbb{R}_k,\mathcal{D}) = rac{1}{N_k} \sum_{x^{(n)} \in \mathbb{R}_k} (y^{(n)} - w_k)^2$$

total cost is the normalized sum over all regions

Cost function

classification cost

find a **decision tree** minimizing the following cost function, this cost function specifies "what is a good decision or regression tree?"

again, calculate the cost per region \mathbb{R}_k



total cost is the normalized sum $ext{cost}(\mathcal{D}) = \sum_k rac{N_k}{N} ext{cost}(\mathbb{R}_k,\mathcal{D})$

Cost function

find a **decision tree** minimizing the following cost function, this cost function specifies "what is a good decision or regression tree?" is sex male?

total cost is the normalized sum $\operatorname{cost}(\mathcal{D}) = \sum_k rac{N_k}{N} \operatorname{cost}(\mathbb{R}_k, \mathcal{D})$

problem

it is sometimes possible to build a tree with **zero cost**:

build a large tree with each instance having its own region (*overfitting*!)



use features such as height, eye color etc, to make perfect prediction on training data

solution

find a decision tree with at most **K tests** minimizing the cost function K tests = K internal node in our binary tree = K+1 leaves (regions)

ves

is age > 9.5?

died

0.05 2%

died

0 17 61%

no

survived

0.73 36%

survived

0.89 2%

is sibsp > 2.5?

Search space

K+1 regions

objective: find a decision tree with K tests minimizing the cost function

the number of full binary trees with K+1 leaves (regions \mathbb{R}_k) is the **Catalan number**



 $\bigwedge \bigwedge \bigwedge \bigwedge \bigwedge \bigwedge \bigwedge$

exponential in K

 $\frac{1}{K+1}\binom{2K}{K}$

we also have a choice of feature x_d for each of K internal node D^K

moreover, for each feature different choices of splitting

bottom line: finding optimal decision tree is an NP-hard combinatorial optimization problem

Greedy heuristic

finding the optimal tree is too difficult, instead use a greedy heuristic to find a good tree recursively split the regions based on a greedy choice of the next test end the recursion if not worth-splitting $x_1 \le t_1$



 $\{\{\mathbb{R}_1,\mathbb{R}_2\},\{\mathbb{R}_3,\{\mathbb{R}_4,\mathbb{R}_5\}\}$

final decision tree in the form of nested list of regions

Choosing tests

the split is greedy because it looks one step ahead this may not lead to the lowest overall cost

```
function greedy-test (\mathbb{R}_{node}, \mathcal{D})
           best-cost = inf
          for each feature d \in \{1,\ldots,D\} and each possible test
                    split \mathbb{R}_{node} into \mathbb{R}_{left}, \mathbb{R}_{right} based on the test
                    \texttt{split-cost} = rac{N_{\mathsf{left}}}{N_{\mathsf{redt}}} \operatorname{cost}(\mathbb{R}_{\mathsf{left}}, \mathcal{D}) + rac{N_{\mathsf{right}}}{N_{\mathsf{redt}}} \operatorname{cost}(\mathbb{R}_{\mathsf{right}}, \mathcal{D})
                    if split-cost < best-cost:
                              best-cost = split-cost
                              \mathbb{R}^*_{\mathsf{left}} = \mathbb{R}_{\mathsf{left}}
                              \mathbb{R}^*_{\mathsf{right}} = \mathbb{R}_{\mathsf{right}}
return \mathbb{R}^*_{\text{left}}, \mathbb{R}^*_{\text{right}}
```

Stopping the recursion

worth-splitting subroutine

if we stop when \mathbb{R}_{node} has zero cost, we may overfit heuristics for stopping the splitting:

- reached a desired depth
- number of examples in $\mathbb{R}_{\mathsf{left}}\;\;\mathsf{or}\;\;\mathbb{R}_{\mathsf{right}}\;\mathsf{is}\;\mathsf{too}\;\mathsf{small}$
- w_k is a good approximation, the cost is small enough
- reduction in cost by splitting is small

$$\mathrm{cost}(\mathbb{R}_{\mathsf{node}},\mathcal{D}) - \left(rac{N_{\mathsf{left}}}{N_{\mathsf{node}}} \mathrm{cost}(\mathbb{R}_{\mathsf{left}},\mathcal{D}) + rac{N_{\mathsf{right}}}{N_{\mathsf{node}}} \mathrm{cost}(\mathbb{R}_{\mathsf{right}},\mathcal{D})
ight)$$





image credit: https://alanjeffares.wordpress.com/tutorials/decision-tree/

revisiting the **classification cost**

ideally we want to optimize the misclassification rate

$$ext{cost}(\mathbb{R}_k,\mathcal{D}) = rac{1}{N_k}\sum_{x^{(n)}\in\mathbb{R}_k}\mathbb{I}(y^{(n)}
eq w_k)$$
)

this may not be the optimal cost for *each step of greedy heuristic*

example both splits have the same misclassification rate (2/8) $\cot(\mathcal{D}) = \frac{1}{N} \sum_{k} \sum_{x^{(n)} \in \mathbb{R}_{k}} \mathbb{I}(y^{(n)} \neq w_{k})$



however the second split *may be* preferable because one region does not need further splitting idea: use a measure for homogeneity of labels in regions

Entropy: a measure of the *unpredictability*

entropy is the **expected amount of information** in observing a random variable

 $H(y) = -\sum_{c=1}^C p(y=c) \log p(y=c)$ averaged on all its possible outcomes

information from two independent events is additive $-\log(p(c)q(d)) = -\log p(c) - \log q(d)$

 $-\log p(y=c)$ is the amount of **information** in observing value c

less probable events are more informative $p(c) < p(c') \Rightarrow -\log p(c) > -\log p(c')$

zero information if p(c)=1

$$I(y = c) = \log(\frac{1}{p(y=c)}) = -log(p(y = c))$$

p(y)

p(y)

a uniform distribution has the highest entropy
$$H(y) = -\sum_{c=1}^{C} \frac{1}{C} \log \frac{1}{C} = \log C$$

a deterministic random variable has the lowest entropy $H(y) = -1 \log(1) = 0$

The most basic concept in the field of Information Theory

note that it is common to use capital letters for random variables (here for consistency we use lower-case)

Mutual information

for two random variables *t*, *y*, their **mutual information** is

the amount of information t conveys about y change in the entropy of y after observing the value of t

how much knowing t reduces uncertainty about y

I(t,y) = H(y) - H(y|t)conditional entropy $\sum_{l=1}^{L} p(t=l)H(x|t=l)$ uncertainty we have in y after seeing t

$$p = \sum_l \sum_c p(y=c,t=l) \log rac{p(y=c,t=l)}{p(y=c)p(t=l)}$$
 this is symmetric wrt y and t

= H(t) - H(t|y) = I(y,t)

mutual information is always positive and zero only if y and t are independent

Classification cost: example

we care about the empirical distribution of labels in each region

$$p_k(y=c) = rac{\sum_{x^{(n)} \in \mathbb{R}_k} \mathbb{I}(y^{(n)}=c)}{N_k}$$

 $\begin{array}{l} \textbf{misclassification cost} \ \ \text{cost}(\mathbb{R}_k,\mathcal{D}) = \frac{1}{N_k}\sum_{x^{(n)}\in\mathbb{R}_k}\mathbb{I}(y^{(n)}\neq w_k) = 1-p_k(w_k) \\ & \quad \text{the most probable class} \ \ w_k = \arg\max_c p_k(c) \end{array}$

$$(.5, 100\%) \bigoplus \bigoplus \bigoplus \mathbb{R}_{node}$$

$$\mathbb{R}_{left} \bigoplus \bigoplus (.25, 50\%) \quad (.75, 50\%) \bigoplus \bigoplus \mathbb{R}_{right}$$

$$p(y = +) = \frac{1}{4}, p(y = -) = \frac{3}{4} \qquad p(y = +) = \frac{3}{4}, p(y = -) = \frac{1}{4}$$

$$w_k = - \qquad \qquad w_k = +$$
misclassification cost = $\frac{4}{8} \cdot \frac{1}{4} + \frac{4}{8} \cdot \frac{1}{4} = \frac{1}{4}$

 $\mathrm{cost}(\mathcal{D}) = \sum_k rac{N_k}{N} \mathrm{cost}(\mathbb{R}_k, \mathcal{D})$

Entropy for classification cost: example

we care about the empirical distribution of labels in each region $p_k(y=c)=rac{\sum_{x^{(n)}\in\mathbb{R}_k}\mathbb{I}(y^{(n)}=c)}{N_k}$

$$\begin{array}{l} \textbf{misclassification cost} \ \ \text{cost}(\mathbb{R}_k,\mathcal{D}) = \frac{1}{N_k}\sum_{x^{(n)}\in\mathbb{R}_k}\mathbb{I}(y^{(n)}\neq w_k) = 1-p_k(w_k) \\ & \quad \text{the most probable class} \ \ w_k = \arg\max_c p_k(c) \end{array}$$

 $entropy \ cost$ $\cost(\mathbb{R}_k,\mathcal{D})=H(y)=-\sum_{c=1}^C p(y=c)\log p(y=c)$ choose the split with the lowest entropy



Entropy for classification cost: example



Entropy for classification cost

we care about the empirical distribution of labels in each region

$$p_k(y=c) = rac{\sum_{x^{(n)} \in \mathbb{R}_k} \mathbb{I}(y^{(n)}=c)}{N_k}$$

 $ext{entropy cost}$ $ext{cost}(\mathbb{R}_k,\mathcal{D})=H(y)$ choose the split with the lowest entropy

change in the cost becomes the mutual information between the test and labels

$$egin{aligned} ext{cost}(\mathbb{R}_{\mathsf{node}},\mathcal{D}) &- \left(rac{N_{\mathsf{left}}}{N_{\mathsf{node}}} \operatorname{cost}(\mathbb{R}_{\mathsf{left}},\mathcal{D}) + rac{N_{\mathsf{left}}}{N_{\mathsf{node}}} \operatorname{cost}(\mathbb{R}_{\mathsf{right}},\mathcal{D})
ight) & I(t,y) = H(y) - H(y|t) \ &\sum_{l=1}^{L} p(t=l) H(y|t=l) \ &= H(y) - \left(p(x_d \geq t) H(y|x_d \geq t) + p(x_d < t) H(y|x_d < t)
ight) = I(y,x > t) \end{aligned}$$

this means by using entropy as our cost, we are choosing the test which is **maximally informative** about labels



Gini index another cost for selecting the *test* in classification

misclassification (error) rate

$$ext{cost}(\mathbb{R}_k,\mathcal{D}) = rac{1}{N_k}\sum_{x^{(n)}\in\mathbb{R}_k}\mathbb{I}(y^{(n)}
eq w_k) = 1-p(w_k)$$

 $\mathsf{entropy} \mid \mathsf{cost}(\mathbb{R}_k,\mathcal{D}) = H(y)$

Gini index

it is the expected error rate

 $ext{cost}(\mathbb{R}_k,\mathcal{D}) = \sum_{c=1}^C p(c)(1-p(c))$ probability of class c probability of error

$$=\sum_{c=1}^{C} p(c) - \sum_{c=1}^{C} p(c)^2 = 1 - \sum_{c=1}^{C} p(c)^2$$



Decision tree example decision tree for Iris dataset



2.5

4.5

5

decision boundaries suggest overfitting

confirmed using a validation set

training accuracy ~ 85% validation accuracy ~ 70%



6.5

6

7

7.5

Decision tree: overfitting

a decision tree can fit any Boolean function (binary classification with binary features)

example: of decision tree representation of a boolean function (D=3)



there are 2^{2^D} such functions, why?

decision tree can perfectly fit our training data

How to solve the problem of overfitting in large decision trees?

idea 1. grow a small tree



problem: substantial reduction in cost may happen after a few steps by stopping early we cannot know this



Decision tree: overfitting & pruning

idea 2. grow a large tree and then prune it

> greedily turn an internal node into a leaf node choice is based on the lowest increase in the cost repeat this until left with the root node pick the best among the above models using a validation set



10

Cross–validatior

Training set

• • • Min + 1 std. err

15

20

26

Best choice



Summary

- model: divide the input into axis-aligned regions
- cost: for regression and classification
- optimization:
 - NP-hard
 - use greedy heuristic
- adjust the cost for the heuristic
 - using entropy (relation to mutual information maximization)
 - using Gini index
- there are variations on decision tree heuristics
 - what we discussed in called *Classification and Regression Trees (CART)*
- Compared to KNN, robust to scaling and noise, fast predictions, more interpretable