

NNC Recap, Decision Trees

Professor Ameet Talwalkar

Slide Credit: Professor Fei Sha

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Decision tree

Registering for Course

- No definitive update, though I am optimistic that we can increase the class size.
- If you are interested in taking the course but are not able to register, please submit HW1 on Tuesday.

Nikos' Office Hours and Section

- He had to slightly change them due to space conflicts.
- New times are Monday 3:30p-4:30p and Wednesday 2:00p-3:00p.
- First Discussion: Friday 12:00pm - 1:50pm, Kinsey Teaching Pavilion 1200B

Homework 1

- Available online
- Octave will be accepted (along with MATLAB)
 - ▶ Programming assignment is due next Tuesday
- Join Piazza if you haven't already

Outline

1 Administration

2 Review of last lecture

- General setup for classification
- Nearest neighbor classifier
- Understanding learning algorithm
- Practical Considerations

3 Decision tree

Multi-class classification

Classify data into one of the multiple categories

- Input (feature vectors): $\mathbf{x} \in \mathbb{R}^D$
- Output (label): $y \in [C] = \{1, 2, \dots, C\}$
- Learning goal: $y = f(\mathbf{x})$

Special case: binary classification

- Number of classes: $C = 2$
- Labels: $\{0, 1\}$ or $\{-1, +1\}$

More terminology

Training data (set)

- N samples/instances: $\mathcal{D}^{\text{TRAIN}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- They are used for learning $f(\cdot)$

Test (evaluation) data

- M samples/instances: $\mathcal{D}^{\text{TEST}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$
- They are used for assessing how well $f(\cdot)$ will do in predicting an unseen $\mathbf{x} \notin \mathcal{D}^{\text{TRAIN}}$

Training data and test data should *not* overlap: $\mathcal{D}^{\text{TRAIN}} \cap \mathcal{D}^{\text{TEST}} = \emptyset$

Algorithm

Nearest neighbor

$$\mathbf{x}(1) = \mathbf{x}_{\text{nn}(\mathbf{x})}$$

where $\text{nn}(\mathbf{x}) \in [\mathbf{N}] = \{1, 2, \dots, \mathbf{N}\}$,

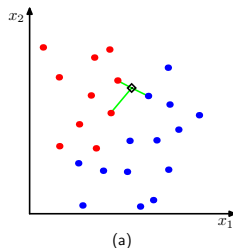
$$\text{nn}(\mathbf{x}) = \arg \min_{n \in [\mathbf{N}]} \|\mathbf{x} - \mathbf{x}_n\|_2^2$$

Classification rule

$$y = f(\mathbf{x}) = y_{\text{nn}(\mathbf{x})}$$

Can extend to KNN classification

- Every neighbor gets a vote;
return the majority vote



Measuring quality of a learning algorithm

- Assume data (\mathbf{x}, y) drawn from *unknown*, joint distribution $p(\mathbf{x}, y)$
- *0/1 loss function* measures mistake on a single data point

$$L(f(\mathbf{x}), y) = \begin{cases} 0 & \text{if } f(\mathbf{x}) = y \\ 1 & \text{if } f(\mathbf{x}) \neq y \end{cases}$$

- Empirical risk (on test set)

$$R_{\mathcal{D}}(f) = \frac{1}{M} \sum_m L(f(\mathbf{x}_m), y_m)$$

- Expected risk

$$R(f) = \mathbb{E}_{(\mathbf{x}, y) \sim p(\mathbf{x}, y)} L(f(\mathbf{x}), y)$$

Bayes binary classifier

It 'cheats' by using the posterior probability $\eta(\mathbf{x}) = p(y = 1|\mathbf{x})$

$$f^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \eta(\mathbf{x}) \geq 1/2 \\ 0 & \text{if } \eta(\mathbf{x}) < 1/2 \end{cases} \quad \text{equivalently } f^*(\mathbf{x}) = \begin{cases} 1 & \text{if } p(y = 1|\mathbf{x}) \geq p(y = 0|\mathbf{x}) \\ 0 & \text{if } p(y = 1|\mathbf{x}) < p(y = 0|\mathbf{x}) \end{cases}$$

Bayes binary classifier

It 'cheats' by using the posterior probability $\eta(\mathbf{x}) = p(y = 1|\mathbf{x})$

$$f^*(\mathbf{x}) = \begin{cases} 1 & \text{if } \eta(\mathbf{x}) \geq 1/2 \\ 0 & \text{if } \eta(\mathbf{x}) < 1/2 \end{cases} \quad \text{equivalently } f^*(\mathbf{x}) = \begin{cases} 1 & \text{if } p(y = 1|\mathbf{x}) \geq p(y = 0|\mathbf{x}) \\ 0 & \text{if } p(y = 1|\mathbf{x}) < p(y = 0|\mathbf{x}) \end{cases}$$

Unsurprisingly, it is optimal

Theorem

For any labeling function $f(\cdot)$, $R(f^) \leq R(f)$.*

Comparing NNC to Bayes optimal classifier

How well does NNC do asymptotically?

Theorem (Cover-Hart inequality)

For the NNC rule f^{NNC} for binary classification, we have,

$$R(f^*) \leq R(f^{\text{NNC}}) \leq 2R(f^*)$$

What does this tell us?

- Shows that as $n \rightarrow \infty$, NNC's expected risk is at worst twice that of the Bayes optimal classifier
- Provides theoretical justification, as NNC is nearly optimal asymptotically

Hypeparameters in NNC

Three practical issues about NNC

- Choosing K , i.e., the number of nearest neighbors (default is 1)
- Choosing the right distance measure (default is Euclidean distance)
- Choosing the scale of each feature since distances depend on units (default is to normalize to zero mean and unit variance)

Those are not specified by the algorithm itself — resolving them requires empirical studies and are task/dataset-specific.

Tuning by using a validation dataset

Training data (set)

- N samples/instances: $\mathcal{D}^{\text{TRAIN}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$
- They are used for learning $f(\cdot)$

Test (evaluation) data

- M samples/instances: $\mathcal{D}^{\text{TEST}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$
- They are used for assessing how well $f(\cdot)$ will do in predicting an unseen $\mathbf{x} \notin \mathcal{D}^{\text{TRAIN}}$

Validation (or development) data

- L samples/instances: $\mathcal{D}^{\text{DEV}} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_L, y_L)\}$
- They are used to optimize hyperparameter(s).

Training data, validation and test data should *not* overlap!

Recipe

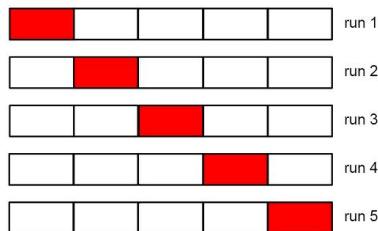
- For each possible value of the hyperparameter (say $K = 1, 3, \dots, 100$)
 - ▶ Train a model using $\mathcal{D}^{\text{TRAIN}}$
 - ▶ Evaluate the performance of the model on \mathcal{D}^{DEV}
- Choose the model with the best performance on \mathcal{D}^{DEV}
- Evaluate the model on $\mathcal{D}^{\text{TEST}}$

Cross-validation

What if we do not have validation data?

- We split the training data into S equal parts.
- We use each part *in turn* as a validation dataset and use the others as a training dataset.
- We choose the hyperparameter such that *on average*, the model performing the best

$S = 5$: 5-fold cross validation



Special case: when $S = N$, this will be leave-one-out.

Recipe

- Split the training data into S equal parts. Denote each part as $\mathcal{D}_s^{\text{TRAIN}}$

Recipe

- Split the training data into S equal parts. Denote each part as $\mathcal{D}_s^{\text{TRAIN}}$
- For each possible value of the hyperparameter (say $K = 1, 3, \dots, 100$)
 - ▶ for every $s \in [1, S]$
 - ★ Train a model using $\mathcal{D}_{\setminus s}^{\text{TRAIN}} = \mathcal{D}^{\text{TRAIN}} - \mathcal{D}_s^{\text{TRAIN}}$
 - ★ Evaluate the performance of the model on $\mathcal{D}_s^{\text{TRAIN}}$
 - ▶ Average the S performance metrics

Recipe

- Split the training data into S equal parts. Denote each part as $\mathcal{D}_s^{\text{TRAIN}}$
- For each possible value of the hyperparameter (say $K = 1, 3, \dots, 100$)
 - ▶ for every $s \in [1, S]$
 - ★ Train a model using $\mathcal{D}_{\setminus s}^{\text{TRAIN}} = \mathcal{D}^{\text{TRAIN}} - \mathcal{D}_s^{\text{TRAIN}}$
 - ★ Evaluate the performance of the model on $\mathcal{D}_s^{\text{TRAIN}}$
 - ▶ Average the S performance metrics
- Choose the hyperparameter corresponding to the best averaged performance

Recipe

- Split the training data into S equal parts. Denote each part as $\mathcal{D}_s^{\text{TRAIN}}$
- For each possible value of the hyperparameter (say $K = 1, 3, \dots, 100$)
 - ▶ for every $s \in [1, S]$
 - ★ Train a model using $\mathcal{D}_{\setminus s}^{\text{TRAIN}} = \mathcal{D}^{\text{TRAIN}} - \mathcal{D}_s^{\text{TRAIN}}$
 - ★ Evaluate the performance of the model on $\mathcal{D}_s^{\text{TRAIN}}$
 - ▶ Average the S performance metrics
- Choose the hyperparameter corresponding to the best averaged performance
- Use the best hyperparameter to train on a model using all $\mathcal{D}^{\text{TRAIN}}$

Recipe

- Split the training data into S equal parts. Denote each part as $\mathcal{D}_s^{\text{TRAIN}}$
- For each possible value of the hyperparameter (say $K = 1, 3, \dots, 100$)
 - ▶ for every $s \in [1, S]$
 - ★ Train a model using $\mathcal{D}_{\setminus s}^{\text{TRAIN}} = \mathcal{D}^{\text{TRAIN}} - \mathcal{D}_s^{\text{TRAIN}}$
 - ★ Evaluate the performance of the model on $\mathcal{D}_s^{\text{TRAIN}}$
 - ▶ Average the S performance metrics
- Choose the hyperparameter corresponding to the best averaged performance
- Use the best hyperparameter to train on a model using all $\mathcal{D}^{\text{TRAIN}}$
- Evaluate the model on $\mathcal{D}^{\text{TEST}}$

Things you need to know

NNC

- Advantages
 - ▶ Computationally, simple and easy to implement – just computing the distance
 - ▶ Theoretically, has good guarantees
- Disadvantages
 - ▶ Computationally intensive for large-scale problems: $O(ND)$ for labeling a data point
 - ▶ We need to “carry” the training data around to perform classification (*nonparametric*).
 - ▶ Choosing the right distance measure, scaling, and K can be involved.

Crucial theoretical concepts loss function, expected risk, empirical risk, Bayes optimal

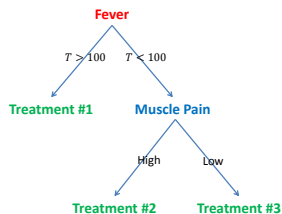
Crucial practical concepts hyperparameters, validation set, cross validation

Outline

- 1 Administration
- 2 Review of last lecture
- 3 Decision tree**
 - Examples
 - Algorithm

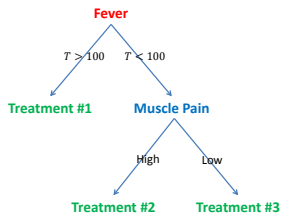
Many decisions are tree structures

Medical treatment

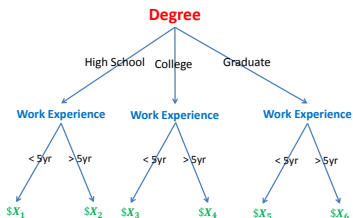


Many decisions are tree structures

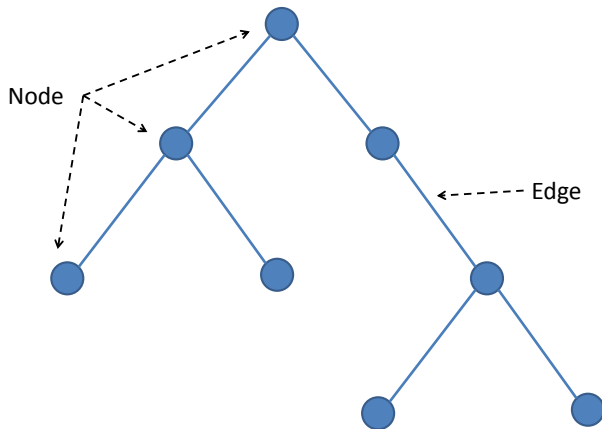
Medical treatment



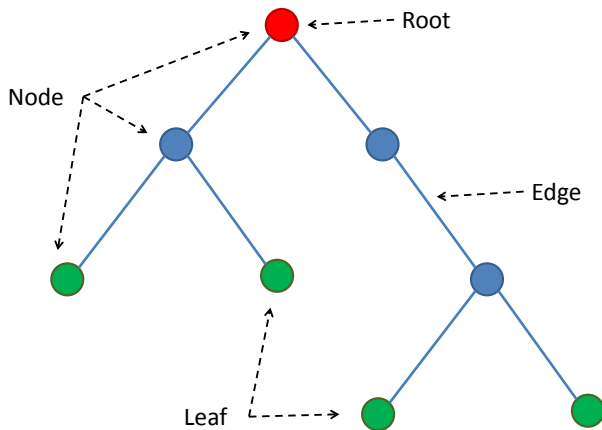
Salary in a company



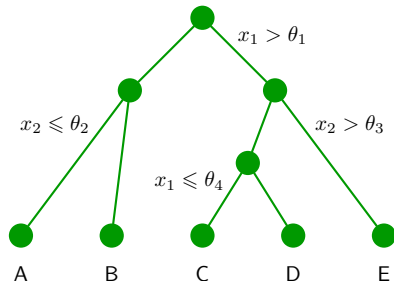
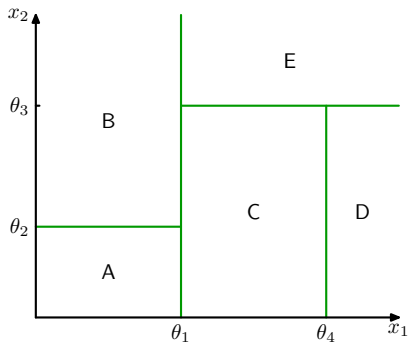
What is a Tree?



Special Names for Nodes in a Tree



A tree partitions the feature space



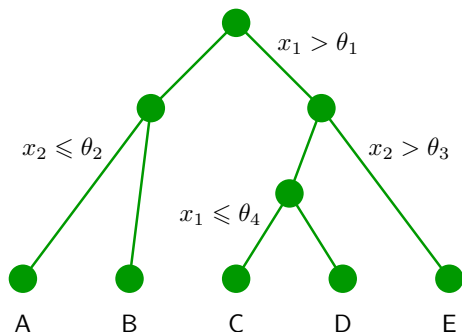
Learning a tree model

Three things to learn:

Learning a tree model

Three things to learn:

- 1 The structure of the tree.
- 2 The threshold values (θ_i).
- 3 The values for the leafs (A, B, \dots).



A tree model for deciding where to eat

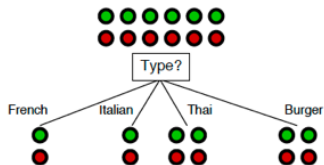
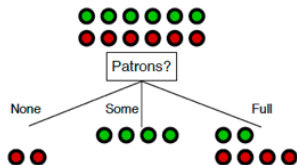
Choosing a restaurant (Example from Russell & Norvig, AIMA)

Example	Attributes										Target
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
X_1	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>0-10</i>	<i>T</i>
X_2	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>30-60</i>	<i>F</i>
X_3	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>Some</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>T</i>
X_4	<i>T</i>	<i>F</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>10-30</i>	<i>T</i>
X_5	<i>T</i>	<i>F</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>>60</i>	<i>F</i>
X_6	<i>F</i>	<i>T</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Italian</i>	<i>0-10</i>	<i>T</i>
X_7	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>F</i>
X_8	<i>F</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Thai</i>	<i>0-10</i>	<i>T</i>
X_9	<i>F</i>	<i>T</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>>60</i>	<i>F</i>
X_{10}	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>Italian</i>	<i>10-30</i>	<i>F</i>
X_{11}	<i>F</i>	<i>F</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>0-10</i>	<i>F</i>
X_{12}	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>30-60</i>	<i>T</i>

Classification of examples is positive (T) or negative (F)

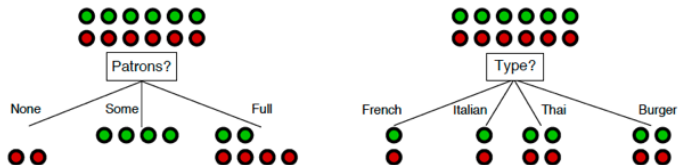
First decision: at the root of the tree

Which attribute to split?



First decision: at the root of the tree

Which attribute to split?



Patrons? is a better choice—gives **information** about the classification

Idea: use information gain to choose
which attribute to split

How to measure information gain?

Idea:


Gaining information reduces uncertainty

Use to entropy to measure uncertainty

If a random variable X has K different values, a_1, a_2, \dots, a_K , its entropy is given by

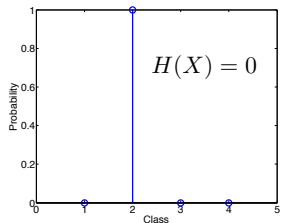
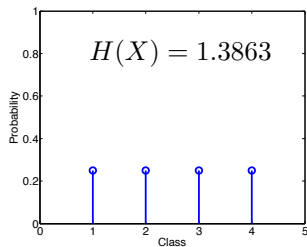
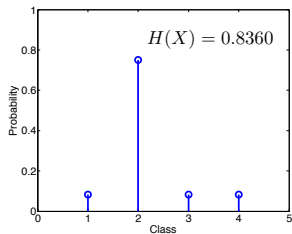
$$H[X] = - \sum_{k=1}^K P(X = a_k) \log P(X = a_k)$$

the base can be 2, though it is not essential (if the base is 2, the unit of the entropy is called "bit")

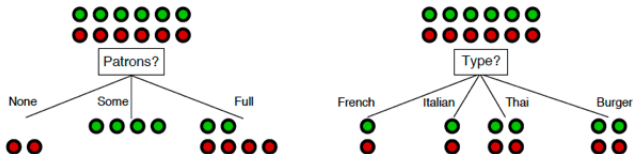


Examples of computing entropy

Entropy



Which attribute to split?



Patrons? is a better choice—gives **information** about the classification

Patron vs. Type?

By choosing Patron, we end up with a partition (3 branches) with smaller entropy, ie, smaller uncertainty (0.45 bit)

By choosing Type, we end up with uncertainty of 1 bit.

Thus, we choose Patron over Type.

Uncertainty if we go with “Patron”

For “None” branch

$$-\left(\frac{0}{0+2} \log \frac{0}{0+2} + \frac{2}{0+2} \log \frac{2}{0+2}\right) = 0$$

For “Some” branch

$$-\left(\frac{4}{4+0} \log \frac{4}{4+0} + \frac{4}{4+0} \log \frac{4}{4+0}\right) = 0$$

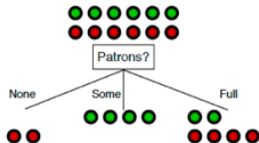
For “Full” branch

$$-\left(\frac{2}{2+4} \log \frac{2}{2+4} + \frac{4}{2+4} \log \frac{4}{2+4}\right) \approx 0.9$$

For choosing “Patrons”

weighted average of each branch: this quantity is called **conditional entropy**

$$\frac{2}{12} * 0 + \frac{4}{12} * 0 + \frac{6}{12} * 0.9 = 0.45$$



Conditional entropy

Definition. Given two random variables X and Y


$$H[Y|X] = \sum_k P(X = a_k)H[Y|X = a_k]$$

In our example

X: the attribute to be split

Y: Wait or not

When $H[Y]$ is fixed, we need only to compare conditional entropy



Relation to information gain

$$\text{GAIN} = H[Y] - H[Y|X]$$

Conditional entropy for Type

For “French” branch

$$-\left(\frac{1}{1+1} \log \frac{1}{1+1} + \frac{1}{1+1} \log \frac{1}{1+1}\right) = 1$$

For “Italian” branch

$$-\left(\frac{1}{1+1} \log \frac{1}{1+1} + \frac{1}{1+1} \log \frac{1}{1+1}\right) = 1$$

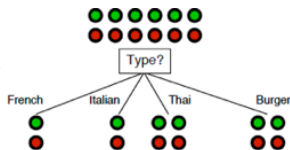
For “Thai” and “Burger” branches

$$-\left(\frac{2}{2+2} \log \frac{2}{2+2} + \frac{2}{2+2} \log \frac{2}{2+2}\right) = 1$$

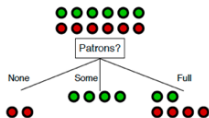
For choosing “Type”

weighted average of each branch:

$$\frac{2}{12} * 1 + \frac{2}{12} * 1 + \frac{4}{12} * 1 + \frac{4}{12} * 1 = 1$$



next split?

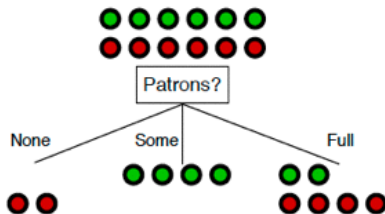


We will look only at the 6 instances with
Patrons == Full

Example	Attributes										
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
X_1	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>0-10</i>	<i>T</i>
X_2	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>30-60</i>	<i>F</i>
X_3	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>Some</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>T</i>
X_4	<i>T</i>	<i>F</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>10-30</i>	<i>T</i>
X_5	<i>T</i>	<i>F</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>>60</i>	<i>F</i>
X_6	<i>F</i>	<i>T</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Italian</i>	<i>0-10</i>	<i>T</i>
X_7	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>F</i>
X_8	<i>F</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Thai</i>	<i>0-10</i>	<i>T</i>
X_9	<i>F</i>	<i>T</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>>60</i>	<i>F</i>
X_{10}	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>Italian</i>	<i>10-30</i>	<i>F</i>
X_{11}	<i>F</i>	<i>F</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>0-10</i>	<i>F</i>
X_{12}	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>30-60</i>	<i>T</i>

Classification of examples is positive (T) or negative (F)

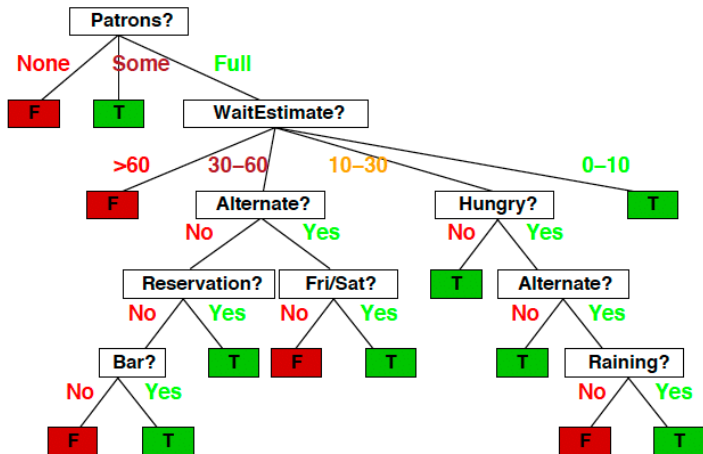
Do we split on “Non” or “Some”?



No, we do not

The decision is deterministic, as seen from the training data

Greedily we build the tree and get this



What is the optimal Tree Depth?

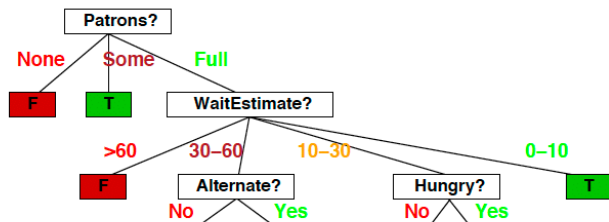
- We need to be careful to pick an appropriate tree depth

What is the optimal Tree Depth?

- We need to be careful to pick an appropriate tree depth
- If the tree is too deep, we can overfit
- If the tree is too shallow, we underfit
- Max depth is a hyperparameter that should be tuned by the data
- Alternative strategy is to create a very deep tree, and then to prune it (see Section 9.2.2 in ESL for details)

Control the size of the tree

We would prune to have a smaller one

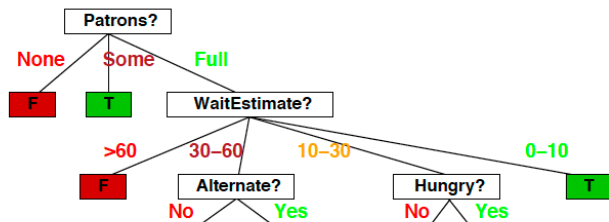


If we stop here, not all training sample would be classified correctly.

More importantly, how do we classify a new instance?

Control the size of the tree

We would prune to have a smaller one



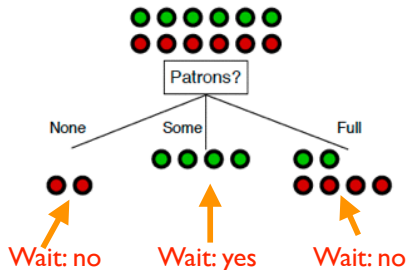
If we stop here, not all training sample would be classified correctly.

More importantly, how do we classify a new instance?

We label the leaves of this smaller tree with **the majority of training samples' labels**

Example

We stop after the root (first node)



Computational Considerations

Numerical Features

- We could split on any feature, with any threshold

Computational Considerations

Numerical Features

- We could split on any feature, with any threshold
- However, for a given feature, the only split points we need to consider are the n values in the training data for this feature.
- If we sort each feature by these n values, we can quickly compute our impurity metric of interest (cross entropy or others)

Computational Considerations

Numerical Features

- We could split on any feature, with any threshold
- However, for a given feature, the only split points we need to consider are the n values in the training data for this feature.
- If we sort each feature by these n values, we can quickly compute our impurity metric of interest (cross entropy or others)
 - ▶ This takes $O(dn \log n)$ time

Computational Considerations

Numerical Features

- We could split on any feature, with any threshold
- However, for a given feature, the only split points we need to consider are the n values in the training data for this feature.
- If we sort each feature by these n values, we can quickly compute our impurity metric of interest (cross entropy or others)
 - ▶ This takes $O(dn \log n)$ time

Categorical Features

- Assuming q distinct categories, there are $2^{q-1} - 1$ possible partitions we can consider.
- However, things simplify in the case of binary classification (or regression), and we can find the optimal split (for cross entropy and Gini) by only considering $q - 1$ possible splits (see Section 9.2.4 in ESL for details).

Summary of learning trees

Advantages of using trees

- Easily interpretable by human (as long as the tree is not too big)
- Computationally efficient
- Handles both numerical and categorical data
- It is parametric thus compact: unlike NNC, we do not have to carry our training instances around
- Building block for various ensemble methods (more on this later)

Disadvantages

- Heuristic training techniques
 - ▶ Finding partition of space that minimizes empirical error is NP-hard
 - ▶ We resort to greedy approaches with limited theoretical underpinnings