

Week 6: Decision Trees and Ensembles

March 13, 2018

1 Decision trees

Imagine you do a newspaper round to help you get through these lean times. On your round, you encounter a number of dogs that either bark or (try to) bite. The dogs are described by the following binary features: *Heavy*, *Smelly*, *Big* and *Growling*. Consider the following set of examples:

Heavy	Smelly	Big	Growling	Bites
No	No	No	No	No
No	No	Yes	No	No
Yes	Yes	No	Yes	No
Yes	No	No	Yes	Yes
No	Yes	Yes	No	Yes
No	No	Yes	Yes	Yes
No	No	No	Yes	Yes
Yes	Yes	No	No	Yes

question 1: What is the entropy of the target value *Bites* in the data?

question 2: Which attribute would the ID3 algorithm choose to use for the root of the tree (without pruning)?

question 3: What is the information gain of the attribute you chose in the previous question?

question 4: Draw the full decision tree that would be learned for this data using ID3 without pruning.

question 5: Suppose three new dogs appear in your round as listed in the table below. Classify them using the decision tree from the previous question.

Dog	Heavy	Smelly	Big	Growling	Bites
Buster	Yes	Yes	Yes	Yes	?
Pluto	No	Yes	No	Yes	?
Zeus	Yes	Yes	No	No	?

question 6: Someone proposes a new scheme to prevent overfitting: she suggests to set a pre-defined maximum depth for the decision trees. When the standard algorithm reaches this depth, it terminates. Could this help to prevent overfitting? Why (not)?

question 7: In the maximum depth scheme introduced above, how would you determine a good value for the maximum depth for a given data set?

question 8: Why can't we apply L1-regularization to this decision tree learning problem?

2 Variational autoencoders

The maximum likelihood principle tells us to optimize the quantity $p(\mathbf{x} \mid \theta)$ as a function of θ (the model parameters).

For complex models, this does not usually lead to a closed form solution and applying gradient descent to θ directly tends to lead to mode collapse. Instead we rely on the following equality:

$$\ln p(\mathbf{x} \mid \theta) = L(\mathbf{q}, \theta) + \text{KL}(\mathbf{q}, \mathbf{p})$$

with

$\mathbf{q}(\mathbf{z})$ any distribution on \mathbf{z}

$$\mathbf{p} = p(\mathbf{z} \mid \mathbf{x}, \theta)$$

$\text{KL}(\mathbf{q}, \mathbf{p})$ the Kullback-Leibler divergence

$$L(\mathbf{q}, \theta) = \mathbb{E}_{\mathbf{q}} \ln \frac{p(\mathbf{x}, \mathbf{z} \mid \theta)}{\mathbf{q}(\mathbf{z})}$$

We will first prove that this equality holds. We start with the right hand side, fill in the components, and derive the left-hand side.

question 9: Fill in the blanks. We have written everything in terms of *expectations* \mathbb{E} to simplify the notation. The expectation is over the random variable \mathbf{z} , while \mathbf{x} has some definite value. Note that $\mathbb{E}f(\mathbf{z}) + \mathbb{E}g(\mathbf{z}) = \mathbb{E}[f(\mathbf{z}) + g(\mathbf{z})]$.

$$\begin{aligned}
L(\mathbf{q}, \theta) + \text{KL}(\mathbf{q}, \mathbf{p}) &= \dots \\
&= \mathbb{E}_{\mathbf{q}} \ln \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{\mathbf{q}(\mathbf{z})} - \mathbb{E}_{\mathbf{q}} \ln \frac{p(\mathbf{z} | \mathbf{x}, \theta)}{\mathbf{q}(\mathbf{z})} \\
&= \dots \\
&= \mathbb{E}_{\mathbf{q}} \ln p(\mathbf{x}, \mathbf{z} | \theta) - \mathbb{E}_{\mathbf{q}} \ln p(\mathbf{z} | \mathbf{x}, \theta) \\
&= \mathbb{E}_{\mathbf{q}} \ln \frac{p(\mathbf{x}, \mathbf{z} | \theta)}{p(\mathbf{z} | \mathbf{x}, \theta)} = \dots \\
&= \mathbb{E}_{\mathbf{q}} \ln p(\mathbf{x} | \theta) = \ln p(\mathbf{x} | \theta)
\end{aligned}$$

In the EM algorithm, we search by alternately optimizing $L(\mathbf{q} | \theta)$ with respect to θ and setting \mathbf{q} equal to \mathbf{p} (so that the KL term becomes zero).

question 10: For the variational autoencoder, we cannot (easily) perform this last step. Why not?

Instead, we *approximate* $p(\mathbf{z} | \mathbf{x}, \theta)$ with a neural network $\mathbf{q}(\mathbf{z} | \mathbf{x})$ that produces a distribution on \mathbf{z} given some \mathbf{x} . This gives us an auto-encoder-like structure. An input is mapped to a distribution $\mathbf{q}(\mathbf{z} | \mathbf{x})$ by the *encoder*. We sample a single \mathbf{z} from this distribution and pass it through the *decoder* $p(\mathbf{x} | \mathbf{z}, \theta)$ to produce a distribution on \mathbf{x} (see the *slides* for diagrams).

To find a way to train such an architecture, we turn again to our decomposition of the likelihood:

$$\ln p(\mathbf{x} | \theta) = L(\mathbf{q}, \theta) + \text{KL}(\mathbf{q}, \mathbf{p}) .$$

Since the KL divergence is always positive, this tells us that

$$\ln p(\mathbf{x} | \theta) \geq L(\mathbf{q}, \theta)$$

for *any* \mathbf{q} we choose. This is why L is called the *variational lower bound*.¹ If we choose our parameters θ and our function \mathbf{q} to maximize L , we are also, indirectly, maximizing $p(\mathbf{x} | \theta)$.²

To do so, we rewrite it into two separate terms: a KL divergence and an expectation:

¹The word *variational* comes from the fact that one of its arguments, \mathbf{q} , is a function (the calculus of functions is called *variational* calculus). For our purposes, this distinction doesn't matter much, since the function \mathbf{q} is defined by a set of parameters, so ultimately we will take the derivative over those parameters, as we are used to.

²How close the lower bound L comes to the true value $p(\mathbf{x} | \theta)$ depends on how well our encoder network \mathbf{q} approximates the true conditional distribution on \mathbf{z} : $p(\mathbf{z} | \mathbf{x}, \theta)$.

$$L(\mathbf{q}, \theta) = -\text{KL}(\mathbf{q}(\mathbf{z} \mid \mathbf{x}), p(\mathbf{z})) + \mathbb{E}_{\mathbf{q}} \ln p(\mathbf{x} \mid \mathbf{z}, \theta)$$

question 11: Show that this equation holds (i.e. rewrite the left part into the right).

Thus, to optimize our variational autoencoder, we should maximize L . In other words, $-L$ is our loss function. The only problem left to solve is that the second term is an expectation (which we cannot compute explicitly).

question 12: How is this solved in practice?