

Week 6: Decision Trees and Variational Autoencoders

<http://mlvu.github.io>

February 29, 2020

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?

$$H(\text{Bites}) = -5/8 \log_2 5/8 - 3/8 \log_2 3/8 \approx 0.9544$$

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

question 3: What is the information gain of the attribute you chose in the previous question? *approximately 0.0487*

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

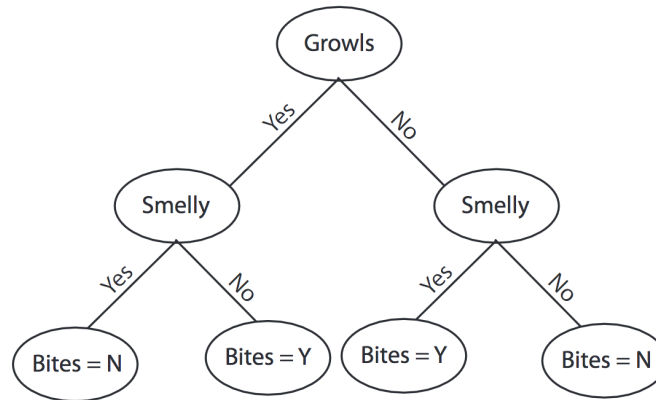


Figure 1: Decision tree

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	No
Pluto	No	Yes	No	Yes	No
Zeus	Yes	Yes	No	No	Yes

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)? It ensures smaller, simpler trees (a smaller model space) and therefore can reduce the risk of overfitting. Overfitting is essentially memorizing too much of your data. Smaller trees can memorize less.

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?

The maximum depth is a *hyperparameter*. We can choose good values for our hyperparameters by splitting our training data into a validation set and trying different values (either by cross validation or just single runs).

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

L1 regularization assumes that your model is described by a real valued vector of parameters (i.e. your model space is continuous.) Decision trees have a *discrete* model space.

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.

Instead, we will rewrite the maximum likelihood objective using the following decomposition.

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

with

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

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

between $\mathbf{q}(z \mid \mathbf{x})$ and $p(z \mid \mathbf{x}, \theta)$

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

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 z , while \mathbf{x} has some definite value. Note that $\mathbb{E}f(z) + \mathbb{E}g(z) = \mathbb{E}[f(z) + g(z)]$.

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

In the EM algorithm, we search by alternately optimizing $L(\mathbf{q} \mid \theta)$ 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? In the variational autoencoder, our model is a neural network that transforms \mathbf{z} into a distribution on \mathbf{x} . To set the KL divergence term equal to zero, we would have to compute $p(\mathbf{z} \mid \mathbf{x}, \theta)$: i.e. a probability distribution on \mathbf{z} that indicates for which \mathbf{z} our observed \mathbf{x} is most likely.

While sampling techniques exist to approximate this kind of distribution, they are costly and can be very inaccurate.

Instead, we approximate $p(\mathbf{z} \mid \mathbf{x}, \theta)$ with a neural network $\mathbf{q}_{\mathbf{v}}(\mathbf{z} \mid \mathbf{x})$ that produces a distribution on \mathbf{z} given some \mathbf{x} . We call the neural network computing $p(\mathbf{x} \mid \mathbf{z}, \theta)$ $\mathbf{p}_{\mathbf{w}}(\mathbf{x} \mid \mathbf{z})$, to make the notation a little more friendly. Here, \mathbf{w} , stands for all parameters of the \mathbf{p} network, and \mathbf{v} stands for all parameters of the \mathbf{q} network.¹

This gives us an auto-encoder-like structure. An input is mapped to a distribution $\mathbf{q}_{\mathbf{v}}(\mathbf{z} \mid \mathbf{x})$ by the **encoder**. We sample a single \mathbf{z} from this distribution and pass it through the **decoder** $\mathbf{p}_{\mathbf{w}}(\mathbf{x} \mid \mathbf{z})$ 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. In our new notation:

$$\ln \mathbf{p}_{\mathbf{w}}(\mathbf{x}) = L(\mathbf{v}, \mathbf{w}) + \text{KL}(\mathbf{q}, \mathbf{p}) .$$

The KL divergence term is difficult to compute: it's an expectation, and it contains the function $\mathbf{p}_{\mathbf{w}}(\mathbf{z} \mid \mathbf{x})$ which requires us to invert the **decoder** neural network (that is, to reason about the inputs given the outputs).

However, because the KL divergence is always positive, we know that

$$\ln \mathbf{p}_{\mathbf{w}}(\mathbf{x}) \geq L(\mathbf{v}, \mathbf{w})$$

for *any* $\mathbf{q}_{\mathbf{v}}$ we choose. This is why L is called the variational *lower bound*.² If we choose our parameters \mathbf{w}, \mathbf{v} to maximize L , we are also, indirectly, max-

¹We've turned θ into \mathbf{w} and added parameters \mathbf{v} for our approximation \mathbf{q} on the conditional distribution on \mathbf{z} . We've also taken the parameters out of the conditional, because we will always talk about the function "given the parameter"; we will never talk about the probability on the parameters themselves.

imizing $\ln p_{\mathbf{w}}(\mathbf{x})$.³

To do so, we rewrite $L(\mathbf{v}, \mathbf{w})$ into two separate terms: a KL divergence and an expectation:

$$L(\mathbf{v}, \mathbf{w}) = -\text{KL}(q_{\mathbf{v}}(\mathbf{z} | \mathbf{x}), p_{\mathbf{v}}(\mathbf{z})) + \mathbb{E}_{q_{\mathbf{v}}} \ln p_{\mathbf{w}}(\mathbf{x} | \mathbf{z})$$

question 11: Show that this equation holds. That is, rewrite the left part into the right. We will assume that all expectations are over q .

$$\begin{aligned} L(\mathbf{v}, \mathbf{w}) &= \mathbb{E}_q \ln \frac{p_{\mathbf{w}}(\mathbf{x}, \mathbf{z})}{q_{\mathbf{v}}(\mathbf{z} | \mathbf{x})} \\ &= \mathbb{E} \ln p_{\mathbf{w}}(\mathbf{x}, \mathbf{z}) - \mathbb{E} \ln q_{\mathbf{v}}(\mathbf{z} | \mathbf{x}) \\ &= \mathbb{E} \ln [p_{\mathbf{w}}(\mathbf{x} | \mathbf{z}) p_{\mathbf{w}}(\mathbf{z})] - \mathbb{E} \ln q_{\mathbf{v}}(\mathbf{z} | \mathbf{x}) \\ &= \mathbb{E} \ln p_{\mathbf{w}}(\mathbf{x} | \mathbf{z}) + \mathbb{E} \ln p_{\mathbf{w}}(\mathbf{z}) - \mathbb{E} \ln q_{\mathbf{v}}(\mathbf{z} | \mathbf{x}) \\ &= \mathbb{E} \ln p_{\mathbf{w}}(\mathbf{x} | \mathbf{z}) - [\mathbb{E} \ln q_{\mathbf{v}}(\mathbf{z} | \mathbf{x}) - \mathbb{E} \ln p_{\mathbf{w}}(\mathbf{z})] \\ &= \mathbb{E} \ln p_{\mathbf{w}}(\mathbf{x} | \mathbf{z}) - \text{KL}(q_{\mathbf{v}}(\mathbf{z} | \mathbf{x}), p_{\mathbf{w}}(\mathbf{z})) \end{aligned}$$

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?

We take a single sample from $q_{\mathbf{v}}(\mathbf{z} | \mathbf{x})$ and use $\ln p_{\mathbf{w}}(\mathbf{x} | \mathbf{z})$ as a (very crude) estimate of the expectation term.

To let the gradient propagate through the sampling, we add a sample from the standard MVN to the input and transform it to a sample from $q_{\mathbf{v}}(\mathbf{z} | \mathbf{x})$ by multiplying by a matrix A (with $\Sigma = AA^T$) and adding the mean.

³The word *variational* comes from the fact that one of its arguments, q , is a function (the calculus of functions is called *variational* calculus). For our purposes, this distinction doesn't matter much, since the function q is defined by a set of parameters \mathbf{v} , so ultimately we will take the derivative over those parameters, as we are used to.

³How close the lower bound L comes to the true value $p_{\mathbf{w}}(\mathbf{x})$ depends on how well our encoder network $q_{\mathbf{v}}$ approximates the true conditional distribution on \mathbf{z} : $p_{\mathbf{w}}(\mathbf{z} | \mathbf{x})$. I.e. how small the KL term in the original decomposition is.