

Data Science from Scratch

by Joel Grus

Copyright © 2015 Joel Grus. All rights reserved.

Printed in the United States of America.

Published by O'Reilly Media, Inc., 1005 Gravenstein Highway North, Sebastopol, CA 95472.

O'Reilly books may be purchased for educational, business, or sales promotional use. Online editions are also available for most titles (<http://safaribooksonline.com>). For more information, contact our corporate/institutional sales department: 800-998-9938 or corporate@oreilly.com.

Editor: Marie Beaugureau

Production Editor: Melanie Yarbrough

Copyeditor: Nan Reinhardt

Proofreader: Eileen Cohen

Indexer: Ellen Troutman-Zaig

Interior Designer: David Futato

Cover Designer: Karen Montgomery

Illustrator: Rebecca Demarest

April 2015: First Edition

Revision History for the First Edition

2015-04-10: First Release

See <http://oreilly.com/catalog/errata.csp?isbn=9781491901427> for release details.

The O'Reilly logo is a registered trademark of O'Reilly Media, Inc. *Data Science from Scratch*, the cover image of a Rock Ptarmigan, and related trade dress are trademarks of O'Reilly Media, Inc.

While the publisher and the author have used good faith efforts to ensure that the information and instructions contained in this work are accurate, the publisher and the author disclaim all responsibility for errors or omissions, including without limitation responsibility for damages resulting from the use of or reliance on this work. Use of the information and instructions contained in this work is at your own risk. If any code samples or other technology this work contains or describes is subject to open source licenses or the intellectual property rights of others, it is your responsibility to ensure that your use thereof complies with such licenses and/or rights.

978-1-491-90142-7

[LSI]

Decision Trees

A tree is an incomprehensible mystery.

—Jim Woodring

DataSciencester’s VP of Talent has interviewed a number of job candidates from the site, with varying degrees of success. He’s collected a data set consisting of several (qualitative) attributes of each candidate, as well as whether that candidate interviewed well or poorly. Could you, he asks, use this data to build a model identifying which candidates will interview well, so that he doesn’t have to waste time conducting interviews?

This seems like a good fit for a *decision tree*, another predictive modeling tool in the data scientist’s kit.

What Is a Decision Tree?

A decision tree uses a tree structure to represent a number of possible *decision paths* and an outcome for each path.

If you have ever played the game **Twenty Questions**, then it turns out you are familiar with decision trees. For example:

- “I am thinking of an animal.”
- “Does it have more than five legs?”
- “No.”
- “Is it delicious?”
- “No.”
- “Does it appear on the back of the Australian five-cent coin?”

- “Yes.”
- “Is it an echidna?”
- “Yes, it is!”

This corresponds to the path:

“Not more than 5 legs” → “Not delicious” → “On the 5-cent coin” → “Echidna!”

in an idiosyncratic (and not very comprehensive) “guess the animal” decision tree (Figure 17-1).

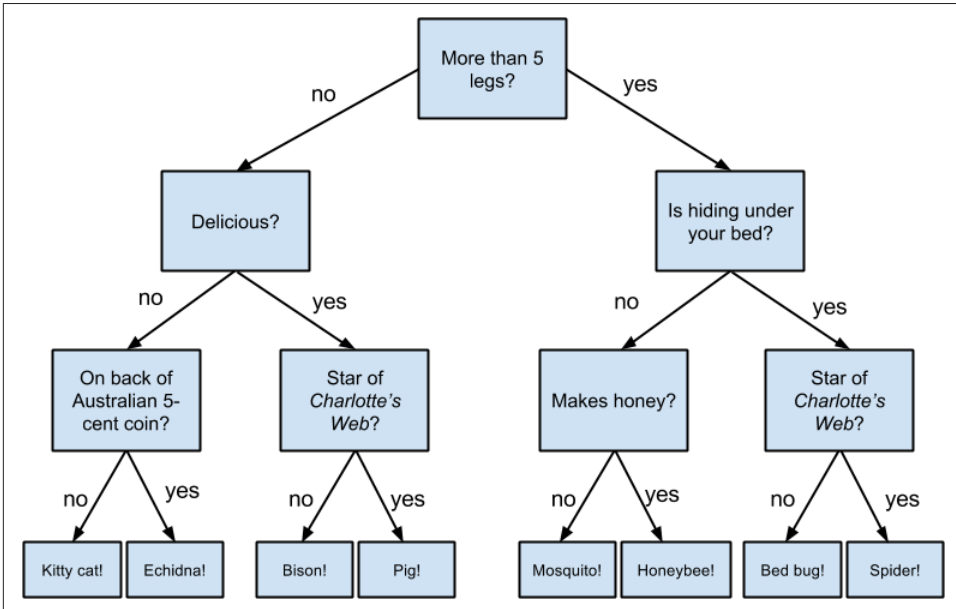


Figure 17-1. A “guess the animal” decision tree

Decision trees have a lot to recommend them. They’re very easy to understand and interpret, and the process by which they reach a prediction is completely transparent. Unlike the other models we’ve looked at so far, decision trees can easily handle a mix of numeric (e.g., number of legs) and categorical (e.g., delicious/not delicious) attributes and can even classify data for which attributes are missing.

At the same time, finding an “optimal” decision tree for a set of training data is computationally a very hard problem. (We will get around this by trying to build a good-enough tree rather than an optimal one, although for large data sets this can still be a lot of work.) More important, it is very easy (and very bad) to build decision trees that are *overfitted* to the training data, and that don’t generalize well to unseen data. We’ll look at ways to address this.

Most people divide decision trees into *classification trees* (which produce categorical outputs) and *regression trees* (which produce numeric outputs). In this chapter, we'll focus on classification trees, and we'll work through the ID3 algorithm for learning a decision tree from a set of labeled data, which should help us understand how decision trees actually work. To make things simple, we'll restrict ourselves to problems with binary outputs like "should I hire this candidate?" or "should I show this website visitor advertisement A or advertisement B?" or "will eating this food I found in the office fridge make me sick?"

Entropy

In order to build a decision tree, we will need to decide what questions to ask and in what order. At each stage of the tree there are some possibilities we've eliminated and some that we haven't. After learning that an animal doesn't have more than five legs, we've eliminated the possibility that it's a grasshopper. We haven't eliminated the possibility that it's a duck. Every possible question partitions the remaining possibilities according to their answers.

Ideally, we'd like to choose questions whose answers give a lot of information about what our tree should predict. If there's a single yes/no question for which "yes" answers always correspond to True outputs and "no" answers to False outputs (or vice versa), this would be an awesome question to pick. Conversely, a yes/no question for which neither answer gives you much new information about what the prediction should be is probably not a good choice.

We capture this notion of "how much information" with *entropy*. You have probably heard this used to mean disorder. We use it to represent the uncertainty associated with data.

Imagine that we have a set S of data, each member of which is labeled as belonging to one of a finite number of classes C_1, \dots, C_n . If all the data points belong to a single class, then there is no real uncertainty, which means we'd like there to be low entropy. If the data points are evenly spread across the classes, there is a lot of uncertainty and we'd like there to be high entropy.

In math terms, if p_i is the proportion of data labeled as class c_i , we define the entropy as:

$$H(S) = -p_1 \log_2 p_1 - \dots - p_n \log_2 p_n$$

with the (standard) convention that $0 \log 0 = 0$.

Without worrying too much about the grisly details, each term $-p_i \log_2 p_i$ is non-negative and is close to zero precisely when p_i is either close to zero or close to one (Figure 17-2).

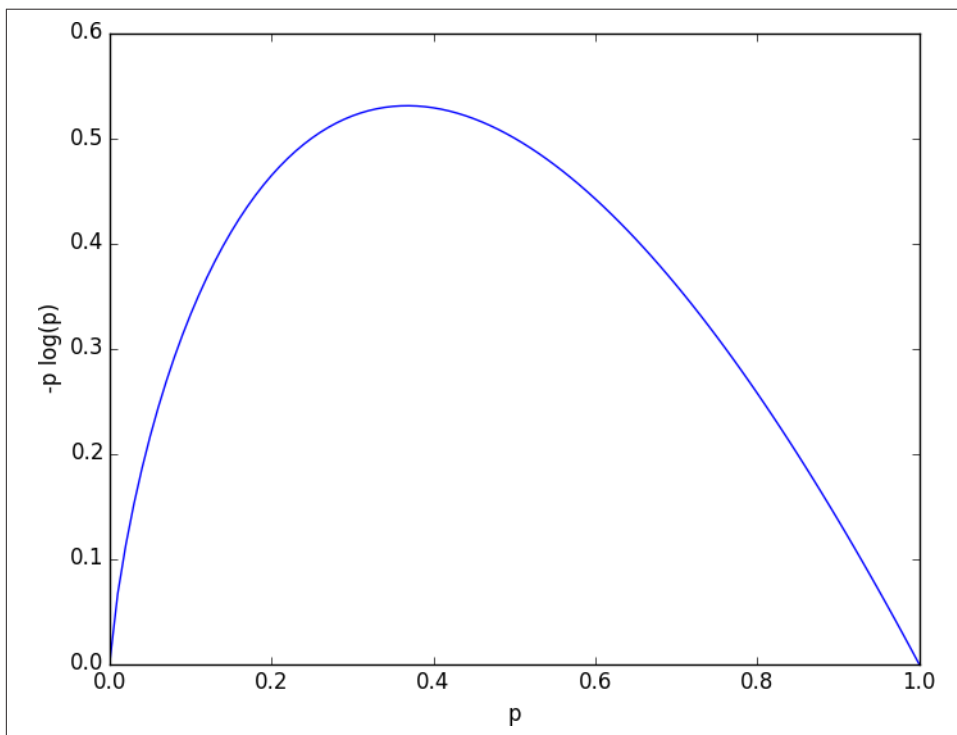


Figure 17-2. A graph of $-p \log p$

This means the entropy will be small when every p_i is close to 0 or 1 (i.e., when most of the data is in a single class), and it will be larger when many of the p_i 's are not close to 0 (i.e., when the data is spread across multiple classes). This is exactly the behavior we desire.

It is easy enough to roll all of this into a function:

```
def entropy(class_probabilities):  
    """given a list of class probabilities, compute the entropy"""  
    return sum(-p * math.log(p, 2)  
              for p in class_probabilities  
              if p)                    # ignore zero probabilities
```

Our data will consist of pairs (input, label), which means that we'll need to compute the class probabilities ourselves. Observe that we don't actually care which label is associated with each probability, only what the probabilities are:

```
def class_probabilities(labels):
    total_count = len(labels)
    return [count / total_count
            for count in Counter(labels).values()]

def data_entropy(labeled_data):
    labels = [label for _, label in labeled_data]
    probabilities = class_probabilities(labels)
    return entropy(probabilities)
```

The Entropy of a Partition

What we’ve done so far is compute the entropy (think “uncertainty”) of a single set of labeled data. Now, each stage of a decision tree involves asking a question whose answer partitions data into one or (hopefully) more subsets. For instance, our “does it have more than five legs?” question partitions animals into those who have more than five legs (e.g., spiders) and those that don’t (e.g., echidnas).

Correspondingly, we’d like some notion of the entropy that results from partitioning a set of data in a certain way. We want a partition to have low entropy if it splits the data into subsets that themselves have low entropy (i.e., are highly certain), and high entropy if it contains subsets that (are large and) have high entropy (i.e., are highly uncertain).

For example, my “Australian five-cent coin” question was pretty dumb (albeit pretty lucky!), as it partitioned the remaining animals at that point into $S_1 = \{\text{echidna}\}$ and $S_2 = \{\text{everything else}\}$, where S_2 is both large and high-entropy. (S_1 has no entropy but it represents a small fraction of the remaining “classes.”)

Mathematically, if we partition our data S into subsets S_1, \dots, S_m containing proportions q_1, \dots, q_m of the data, then we compute the entropy of the partition as a weighted sum:

$$H = q_1 H(S_1) + \dots + q_m H(S_m)$$

which we can implement as:

```
def partition_entropy(subsets):
    """find the entropy from this partition of data into subsets
    subsets is a list of lists of labeled data"""

    total_count = sum(len(subset) for subset in subsets)

    return sum( data_entropy(subset) * len(subset) / total_count
               for subset in subsets )
```



One problem with this approach is that partitioning by an attribute with many different values will result in a very low entropy due to overfitting. For example, imagine you work for a bank and are trying to build a decision tree to predict which of your customers are likely to default on their mortgages, using some historical data as your training set. Imagine further that the data set contains each customer's Social Security number. Partitioning on SSN will produce one-person subsets, each of which necessarily has zero entropy. But a model that relies on SSN is *certain* not to generalize beyond the training set. For this reason, you should probably try to avoid (or bucket, if appropriate) attributes with large numbers of possible values when creating decision trees.

Creating a Decision Tree

The VP provides you with the interviewee data, consisting of (per your specification) pairs (input, label), where each input is a dict of candidate attributes, and each label is either True (the candidate interviewed well) or False (the candidate interviewed poorly). In particular, you are provided with each candidate's level, her preferred language, whether she is active on Twitter, and whether she has a PhD:

```
inputs = [
    ({'level': 'Senior', 'lang': 'Java', 'tweets': 'no', 'phd': 'no'}, False),
    ({'level': 'Senior', 'lang': 'Java', 'tweets': 'no', 'phd': 'yes'}, False),
    ({'level': 'Mid', 'lang': 'Python', 'tweets': 'no', 'phd': 'no'}, True),
    ({'level': 'Junior', 'lang': 'Python', 'tweets': 'no', 'phd': 'no'}, True),
    ({'level': 'Junior', 'lang': 'R', 'tweets': 'yes', 'phd': 'no'}, True),
    ({'level': 'Junior', 'lang': 'R', 'tweets': 'yes', 'phd': 'yes'}, False),
    ({'level': 'Mid', 'lang': 'R', 'tweets': 'yes', 'phd': 'yes'}, True),
    ({'level': 'Senior', 'lang': 'Python', 'tweets': 'no', 'phd': 'no'}, False),
    ({'level': 'Senior', 'lang': 'R', 'tweets': 'yes', 'phd': 'no'}, True),
    ({'level': 'Junior', 'lang': 'Python', 'tweets': 'yes', 'phd': 'no'}, True),
    ({'level': 'Senior', 'lang': 'Python', 'tweets': 'yes', 'phd': 'yes'}, True),
    ({'level': 'Mid', 'lang': 'Python', 'tweets': 'no', 'phd': 'yes'}, True),
    ({'level': 'Mid', 'lang': 'Java', 'tweets': 'yes', 'phd': 'no'}, True),
    ({'level': 'Junior', 'lang': 'Python', 'tweets': 'no', 'phd': 'yes'}, False)
]
```

Our tree will consist of *decision nodes* (which ask a question and direct us differently depending on the answer) and *leaf nodes* (which give us a prediction). We will build it using the relatively simple *ID3* algorithm, which operates in the following manner. Let's say we're given some labeled data, and a list of attributes to consider branching on.

- If the data all have the same label, then create a leaf node that predicts that label and then stop.

- If the list of attributes is empty (i.e., there are no more possible questions to ask), then create a leaf node that predicts the most common label and then stop.
- Otherwise, try partitioning the data by each of the attributes
- Choose the partition with the lowest partition entropy
- Add a decision node based on the chosen attribute
- Recur on each partitioned subset using the remaining attributes

This is what's known as a “greedy” algorithm because, at each step, it chooses the most immediately best option. Given a data set, there may be a better tree with a worse-looking first move. If so, this algorithm won't find it. Nonetheless, it is relatively easy to understand and implement, which makes it a good place to begin exploring decision trees.

Let's manually go through these steps on the interviewee data set. The data set has both True and False labels, and we have four attributes we can split on. So our first step will be to find the partition with the least entropy. We'll start by writing a function that does the partitioning:

```
def partition_by(inputs, attribute):
    """each input is a pair (attribute_dict, label).
    returns a dict : attribute_value -> inputs"""
    groups = defaultdict(list)
    for input in inputs:
        key = input[0][attribute] # get the value of the specified attribute
        groups[key].append(input) # then add this input to the correct list
    return groups
```

and one that uses it to compute entropy:

```
def partition_entropy_by(inputs, attribute):
    """computes the entropy corresponding to the given partition"""
    partitions = partition_by(inputs, attribute)
    return partition_entropy(partitions.values())
```

Then we just need to find the minimum-entropy partition for the whole data set:

```
for key in ['level', 'lang', 'tweets', 'phd']:
    print key, partition_entropy_by(inputs, key)

# level 0.693536138896
# lang 0.860131712855
# tweets 0.788450457308
# phd 0.892158928262
```

The lowest entropy comes from splitting on `level`, so we'll need to make a subtree for each possible `level` value. Every `Mid` candidate is labeled `True`, which means that the `Mid` subtree is simply a leaf node predicting `True`. For `Senior` candidates, we have a mix of `Trues` and `Falses`, so we need to split again:


```

senior_inputs = [(input, label)
                  for input, label in inputs if input["level"] == "Senior"]

for key in ['lang', 'tweets', 'phd']:
    print key, partition_entropy_by(senior_inputs, key)

# lang 0.4
# tweets 0.0
# phd 0.950977500433

```

This shows us that our next split should be on tweets, which results in a zero-entropy partition. For these Senior-level candidates, “yes” tweets always result in True while “no” tweets always result in False.

Finally, if we do the same thing for the Junior candidates, we end up splitting on phd, after which we find that no PhD always results in True and PhD always results in False.

Figure 17-3 shows the complete decision tree.

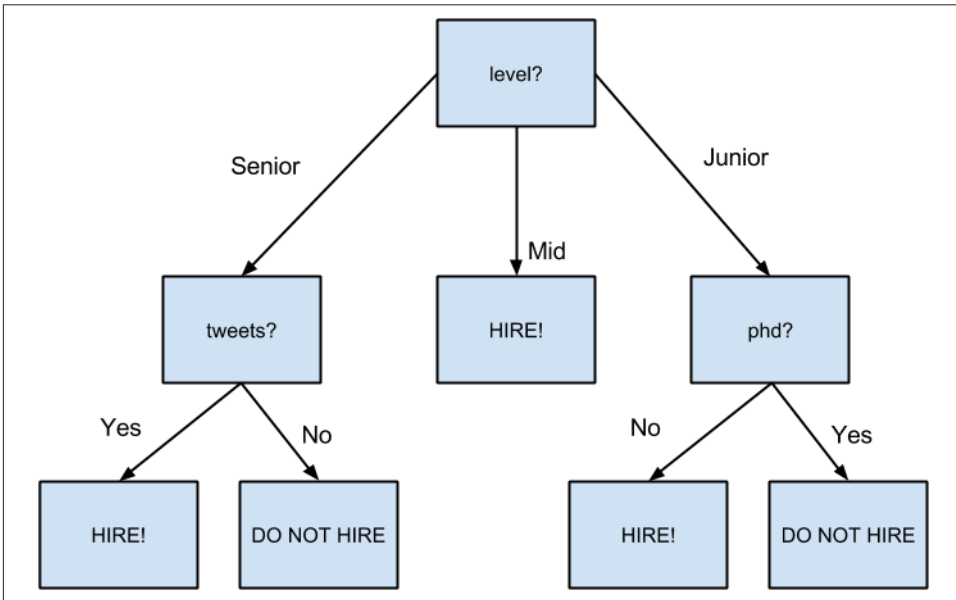


Figure 17-3. The decision tree for hiring

Putting It All Together

Now that we’ve seen how the algorithm works, we would like to implement it more generally. This means we need to decide how we want to represent trees. We’ll use pretty much the most lightweight representation possible. We define a *tree* to be one of the following:

- True
- False
- a tuple (attribute, subtree_dict)

Here True represents a leaf node that returns True for any input, False represents a leaf node that returns False for any input, and a tuple represents a decision node that, for any input, finds its attribute value, and classifies the input using the corresponding subtree.

With this representation, our hiring tree would look like:

```
('level',
 {'Junior': ('phd', {'no': True, 'yes': False}),
  'Mid': True,
  'Senior': ('tweets', {'no': False, 'yes': True})})
```

There’s still the question of what to do if we encounter an unexpected (or missing) attribute value. What should our hiring tree do if it encounters a candidate whose level is “Intern”? We’ll handle this case by adding a None key that just predicts the most common label. (Although this would be a bad idea if None is actually a value that appears in the data.)

Given such a representation, we can classify an input with:

```
def classify(tree, input):
    """classify the input using the given decision tree"""

    # if this is a leaf node, return its value
    if tree in [True, False]:
        return tree

    # otherwise this tree consists of an attribute to split on
    # and a dictionary whose keys are values of that attribute
    # and whose values of are subtrees to consider next
    attribute, subtree_dict = tree

    subtree_key = input.get(attribute)    # None if input is missing attribute

    if subtree_key not in subtree_dict:    # if no subtree for key,
        subtree_key = None                # we'll use the None subtree

    subtree = subtree_dict[subtree_key]   # choose the appropriate subtree
    return classify(subtree, input)      # and use it to classify the input
```

All that’s left is to build the tree representation from our training data:

```
def build_tree_id3(inputs, split_candidates=None):

    # if this is our first pass,
    # all keys of the first input are split candidates
```

```

if split_candidates is None:
    split_candidates = inputs[0][0].keys()

# count Trues and Falses in the inputs
num_inputs = len(inputs)
num_trues = len([label for item, label in inputs if label])
num_falses = num_inputs - num_trues

if num_trues == 0: return False    # no Trues? return a "False" leaf
if num_falses == 0: return True    # no Falses? return a "True" leaf

if not split_candidates:          # if no split candidates left
    return num_trues >= num_falses # return the majority leaf

# otherwise, split on the best attribute
best_attribute = min(split_candidates,
                     key=partial(partition_entropy_by, inputs))

partitions = partition_by(inputs, best_attribute)
new_candidates = [a for a in split_candidates
                  if a != best_attribute]

# recursively build the subtrees
subtrees = { attribute_value : build_tree_id3(subset, new_candidates)
             for attribute_value, subset in partitions.iteritems() }

subtrees[None] = num_trues > num_falses    # default case

return (best_attribute, subtrees)

```

In the tree we built, every leaf consisted entirely of True inputs or entirely of False inputs. This means that the tree predicts perfectly on the training data set. But we can also apply it to new data that wasn't in the training set:

```

tree = build_tree_id3(inputs)

classify(tree, { "level" : "Junior",
                 "lang" : "Java",
                 "tweets" : "yes",
                 "phd" : "no" } )    # True

classify(tree, { "level" : "Junior",
                 "lang" : "Java",
                 "tweets" : "yes",
                 "phd" : "yes" } )    # False

```

And also to data with missing or unexpected values:

```

classify(tree, { "level" : "Intern" } ) # True
classify(tree, { "level" : "Senior" } ) # False

```



Since our goal was mainly to demonstrate *how* to build a tree, we built the tree using the entire data set. As always, if we were really trying to create a good model for something, we would have (collected more data and) split the data into train/validation/test sub-sets.

Random Forests

Given how closely decision trees can fit themselves to their training data, it's not surprising that they have a tendency to overfit. One way of avoiding this is a technique called *random forests*, in which we build multiple decision trees and let them vote on how to classify inputs:

```
def forest_classify(trees, input):
    votes = [classify(tree, input) for tree in trees]
    vote_counts = Counter(votes)
    return vote_counts.most_common(1)[0][0]
```

Our tree-building process was deterministic, so how do we get random trees?

One piece involves bootstrapping data (recall “**Digression: The Bootstrap**” on page 183). Rather than training each tree on all the inputs in the training set, we train each tree on the result of `bootstrap_sample(inputs)`. Since each tree is built using different data, each tree will be different from every other tree. (A side benefit is that it's totally fair to use the nonsampled data to test each tree, which means you can get away with using all of your data as the training set if you are clever in how you measure performance.) This technique is known as *bootstrap aggregating* or *bagging*.

A second source of randomness involves changing the way we chose the `best_attribute` to split on. Rather than looking at all the remaining attributes, we first choose a random subset of them and then split on whichever of those is best:

```
# if there's already few enough split candidates, look at all of them
if len(split_candidates) <= self.num_split_candidates:
    sampled_split_candidates = split_candidates
# otherwise pick a random sample
else:
    sampled_split_candidates = random.sample(split_candidates,
                                             self.num_split_candidates)

# now choose the best attribute only from those candidates
best_attribute = min(sampled_split_candidates,
                     key=partial(partition_entropy_by, inputs))

partitions = partition_by(inputs, best_attribute)
```

This is an example of a broader technique called *ensemble learning* in which we combine several *weak learners* (typically high-bias, low-variance models) in order to produce an overall strong model.

Random forests are one of the most popular and versatile models around.

For Further Exploration

- scikit-learn has many **Decision Tree** models. It also has an **ensemble** module that includes a `RandomForestClassifier` as well as other ensemble methods.
- We barely scratched the surface of decision trees and their algorithms. **Wikipedia** is a good starting point for broader exploration.