

Chapter 7

Naïve Bayes and Decision Trees

[Slides \(pdf\)](#)

[Video \(Panopto\)](#)

In the previous chapter we dove into the workings of the Logistic Regression Model. As you might remember, the Logistic Regression model utilizes the sigmoid function to map out arbitrary input values to output values between 0 and 1. How close an output was to either 0 or 1 gave us a certainty about classifying it, while an output closer to 0.5 meant higher ambiguity. In this chapter we will introduce another way to explore these probabilities.

7.1 Naïve Bayes

Naïve Bayes is an alternative way to compute a class probability for a given input. The formula for a given input x and its output class y is as follows:

Definition 7.1: Bayes Rule

Bayes Theorem computes the probability of a class using the formula below. If we want to compute the probability that y is a positive review given review x , we multiply the probability of x given that y is positive with the probability that y is positive, and divide it with the probability of x alone.

$$P(y = +1|x) = \frac{P(x|y=+1)P(y=+1)}{P(x)}$$

We use this equation to compute how likely a review is positive or negative:

$$\frac{P(\text{"The sushi & everything else was awesome!"} | y = +1) P(y = +1)}{P(\text{"The sushi & everything else was awesome!"})}$$

We may discard the divisor, as we are only comparing which class's probability is greater. Now, think about how you would approach this problem. We would need to plug in a probability for each part of the Bayes Theorem. Do we know $P(\text{"The sushi everything else was awesome!"} | y=+1)$, the probability of this specific review, given that we are told the review is positive? Most likely not, as the sentence "The sushi everything else was awesome!" is unique, and it only appears in one review.

Thus, we make the naïve assumption that every word's probability is independent from each other: Instead of computing this entire sentence's probability, we compute the probability of every word occurring alongside each other. So we compute the probability of "The" given the review is positive AND the probability of "sushi" given the review is positive, and so forth:

$$P(\text{"The sushi everything else was awesome!"} | y=+1) = P(\text{The} | y=+1) * P(\text{sushi} | y=+1) * P(\text{ | } | y=+1) * P(\text{everything} | y=+1) * P(\text{else} | y=+1) * P(\text{was} | y=+1) * P(\text{awesome} | y=+1)$$

Our final model is thus:

$$P(y|x_1, x_2, x_3, \dots, x_d) = \prod_{j=1}^d P(x_j|y)P(y)$$

There are a number of issues with this approach. First, since we are multiplying so many probabilities, our final product will be a very long decimal number. Since this decimal number is so long, we might end up with floating point overflow. We can overcome this by taking the log of each probabilities, such that we compute a sum instead of a product. Another issue that we encounter with using products is that if we encounter a word we haven't seen before, its probability will be 0 and so the entire product will be 0. Laplacian Smoothing, adding a constant to each term to avoid multiplying by 0, can be used in this case.

Let us now compare the two models we have learned for computing classification probabilities thus far:

Logistic Regression:

$$P(y = +1|x, w) = \frac{1}{1 + e^{-w^T h(x_i)}}$$

Naïve Bayes:

$$P(y = +1|x_1, x_2, x_3, \dots, x_d) = \prod_{j=1}^d P(x_j|y = +1)P(y = +1)$$

While the Logistic Regression model is discriminative, the Naïve Bayes model is generative.

Definition 7.2: Discriminative Model

A model is discriminative when it only cares about finding and optimizing a decision boundary.

Definition 7.3: Generative Model

A model is generative when it defines a distribution for generating x .

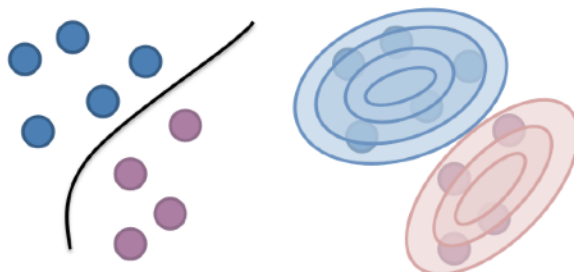


Figure 7.1: On the left visualize a discriminative model, where we only care about optimizing a decision boundary. On the right, we have a generative model, where the likelihood of each class is also captured.

The Naïve Bayes model is considered a generative model because it computes each input's probability given the output $P(x_j|y = +1)$, so it can be used to map out a probability distribution for each x . This is unlike the Logistic Regression model, that does not compute the probability of each input, but rather solely the probability of the output given the input. Usually a discriminative model like Logistic Regression will do better, but generative models are useful for modeling distributions.

7.2 Decision Trees

The next application of machine learning we will explore will be decision trees. Take a look at the decision tree below:

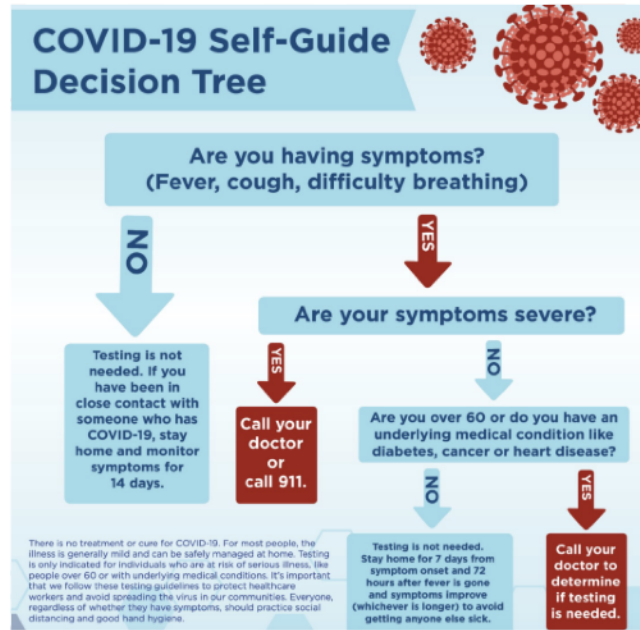


Figure 7.2: A decision tree on COVID-19 safety guidelines based on a series of questions. Source: [Holzer](#)

Definition 7.4: Decision Tree

A decision tree is a series of questions that lead to multiple outcomes such that we can explore each question further.

Decision trees are most applicable for nonlinear decision boundaries that occur along a set of distinct criteria. Let's consider a Loan Application for buying a house as an example. To make an informed decision about a candidate, a loan application may ask about credit history, income, the lease term, and personal information. So, an example dataset might look like this:

Credit	Term	Income	y
excellent	3 yrs	high	safe
fair	5 yrs	low	risky
fair	3 yrs	high	safe
poor	5 yrs	high	risky
excellent	3 yrs	low	safe
fair	5 yrs	low	safe
poor	3 yrs	high	risky
poor	5 yrs	low	safe
fair	3 yrs	high	safe

Figure 7.3: Let our training dataset be n number of inputs and output pairs. Each input has 3 features: credit, term and income and one output, y .

After training on the train dataset, we aim to reach an accurate prediction \hat{y} of whether a new candidate is safe or risky based on their loan application:

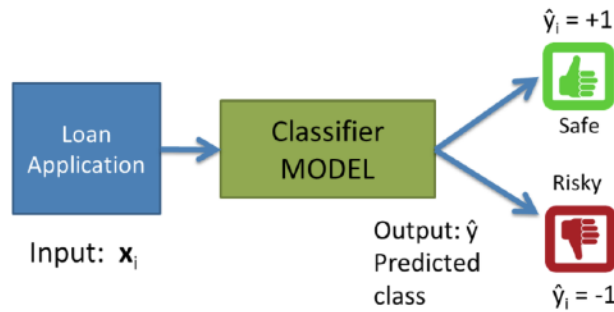


Figure 7.4: A loan application is fed into the model as input x_i . The output will be either $\hat{y}_i = +1$ (safe) or $\hat{y}_i = -1$ (risky).

A decision tree for this example can be set up like so:

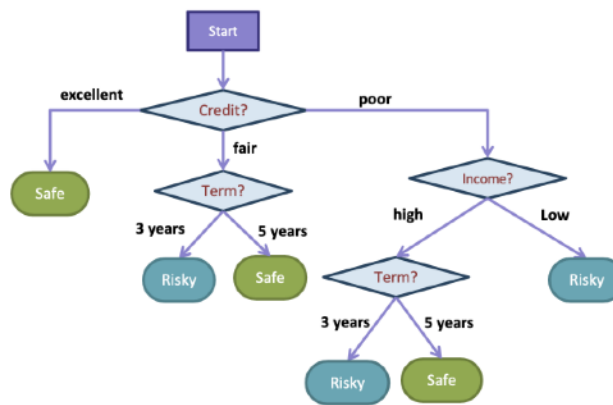


Figure 7.5: A decision tree for determining whether a candidate is safe or risky to lend a loan to. The branch/internal nodes split into possible values of a feature, while the leaf nodes are the final decision, the output class.

Definition 7.5: Tree Terminology

A **Tree** in computer science is a datastructure characterized by an interconnected set of nodes, where no connection between nodes creates a cycle.

Node: any datapoint in the tree.

Root: the first datapoint in the tree.

Leaf: one of the last datapoints in the tree.

In this example, the model first splits up loan applications based on credit history. If the credit history is excellent, a candidate's application is considered safe. However, if their credit history is poor, the tree splits into a more complex structure as more questions are required to determine if the candidate should receive a loan.

At this point, you may be wondering why we use credit history as the first criteria to judge an applicant, why we split into further questions on certain nodes, and how a lease term can determine an application's safety. These are all choices our model made to optimize its tree such that the tree would output the most probable prediction \hat{y} based on its training on the train data. We will now examine how these choices are made in the building of a decision tree.

7.2.1 Building the Decision Tree

How does a model decide on the optimal decision tree structure? Recall that our training data is what our machine learning model uses to learn on. Given a set of x and y training pairs, we will build a decision tree that aims to lead as many train inputs x to their correct train outputs y as would perform the best in the real world. Since we want to avoid overfitting on the train set, this doesn't necessarily mean that we will build our tree such that every training input x will always lead to its exact correct training output y , so we will also have to decide a good stopping point for growing our tree. The end goal is that when new, unseen test data x is fed into the tree, we can achieve a decently accurate prediction \hat{y} .

Let's first start with the very first node in our tree, the root node:

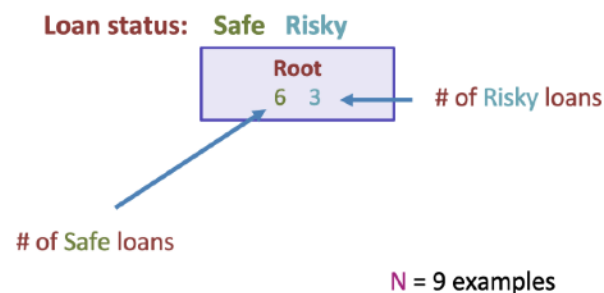


Figure 7.6: The root node of the decision tree stores information about the safe loans and the risky loans. In total, there are 9 loan applications in our dataset.

Since we have just one node, this one node is the root node and it is a leaf node at the same time. In a decision tree, the decisions are made at the leaf nodes. The outputted decision is based on the majority class of the training set inside that node. So if our tree is just a root, then our tree would always output "Safe" since there are 6 safe loans in the training set and 3 risky loans. So, on future unseen test data, we will always output "Safe".

On the train data, this gives us a $6/9 = 66.66\%$ accuracy. If we decide this isn't good enough, we will add a second layer to the tree.

Suppose then, that we decide to split on a candidate's credit history. Credit history can be either "excellent", "fair" or "poor".

Since we split on credit history, we add a node to our tree for each type of credit history. For each of these nodes, we now store how many of the loans are safe and how many are risky. Recall in our original dataset in Figure 7.3 there are 2 safe loans and 0 risky loans that have excellent credit history, 3 safe loans and 1 risky loan with fair credit history, and 1 safe loans and 2 risky loans with poor credit history.

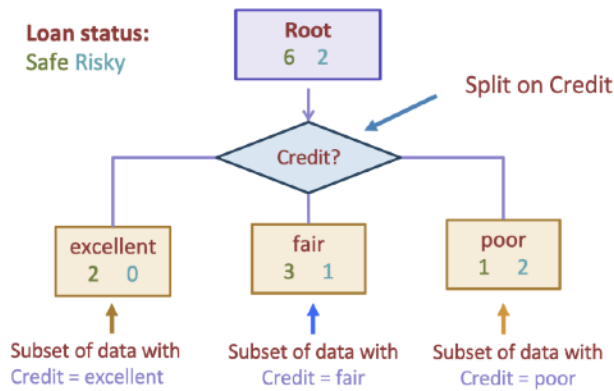


Figure 7.7: The tree splits the training data into different levels of credit history.

If we end our tree here, we must finish it off with leaf nodes that have final output classes in them as the final decisions. Remember that these are the majorities of each node.

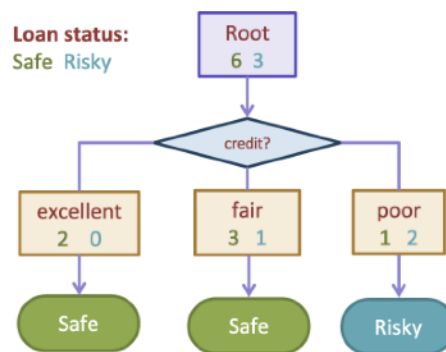


Figure 7.8: The majority class of each node determines the final output.

Now, our training accuracy is $7/9 = 77.77\%$ since 7 out of the 9 loan applications are correctly classified, and 2 are not. What if we were to split on Term first instead of Credit? Let's compare the two choices:

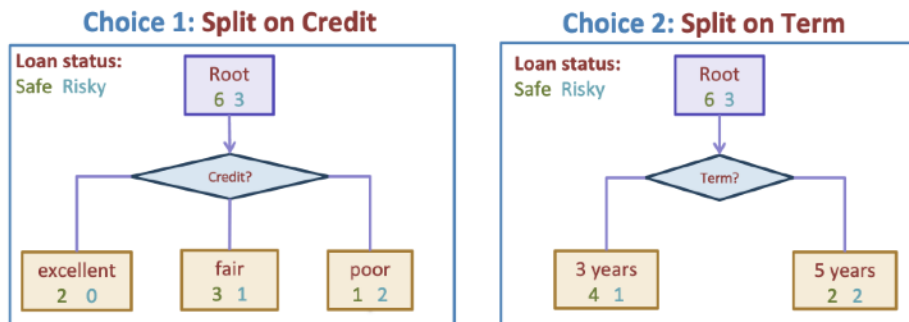


Figure 7.9: A split on Credit versus a split on Term in the decision tree.

In Choice 1 when we split on credit, we have a $7/9 = 77.77\%$ accuracy. Let's calculate the accuracy for Choice 2. For a term of 3 years, we have 4 safe loans and 1 risky loan. Safe loans make up the majority here, so anything with a term of 3 years will be classified as "Safe". For a term of 5 years, there are 2 safe loans and 2 risky loans. Since there is a tie, we can choose either class. Let's designate the output of a 5 year term as a "Risky" loan. Since 4 loan applications for 3 year terms are classified as "Safe" when they actually are and 2 loan applications for 5 year terms are classified as "Risky" when they actually are, $4 + 2 = 6$ loan applications are correctly classified. $6/9 = 66.66\%$ of loan applications are classified correctly. Likewise, we can also use an error metric here: 33.33% of loan applications are classified incorrectly if we split on the Term feature, and 22.22% of loan applications are classified incorrectly if we split on the Credit feature. Since the Credit feature split yields a lower error rate, our tree will split on Credit first. This decision making is how the model splits a node in a tree and it can be summarized by the following pseudocode:

Split(node)

- Given a subset of data M in node
- For each feature h_i :
 - Compute classification error for a split of M according to feature h_i :
- Chose feature $h^*(x)$ with lowest classification error and expand the tree to include the children of current node after the split

Figure 7.10: The algorithm for selecting the best feature a node should be split on.

Our model will choose between features to split on for every node in the tree, until the number of datapoints in each node or the error rate reaches a certain threshold to stop building the tree. This algorithm can be summarized as:

BuildTree(node)

- If the number of datapoints at the current node or the classification error is within a certain threshold:
 - Stop
- Else:
 - Split(node)
 - For child in node:
 - BuildTree(child)

Figure 7.11: The algorithm for building the tree.

Notice that the Decision Tree algorithm is greedy: It aims to optimize the classification error at each node. As a result, the final result won't be globally optimal, but it guarantees computational efficiency. Also take note that the Decision Tree algorithm is recursive: From the current node, if we decide to further expand the tree, we repeat the same operations in the child node.

How do we decide when to stop? Going back to Figure 7.7, we observe that for applicants with excellent credit history, all two of them would be safe loans and none of them would be risky. So, the majority class classifier would suffice here, because it would yield no error. But when the Credit is fair or poor, we had some safe loans and some risky loans. As there is no absolute majority for either of these nodes, we could recursively treat them as roots of their own tree which we continue to build below. Thus by further splitting these nodes, we grow a more complex and precise path in our tree with the hopes of lowering the error rate.

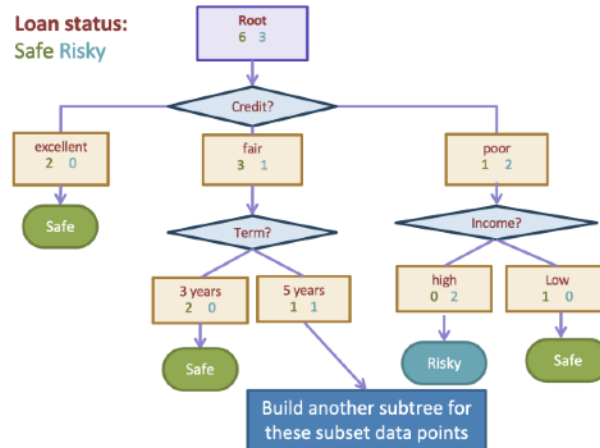


Figure 7.12: Every subtree build from each node will be a subset of the original trainset where the condition of the node is true for the remainder of the tree.

We can set a certain error threshold and stop growing our tree until every leaf yields that error or less.

7.2.2 Numeric v.s. Categorical Features

In our Loan Application example, we have been looking at data without numeric value. For example, a candidate's credit history can be of the values "good", "fair", or "poor". There are 3 main datatypes a machine learning algorithm can learn from.

Definition 7.6: Numeric

Numeric data is data that has numerical value, such as square footage of a house.

Definition 7.7: Ordinal

Ordinal data is data that has ordered categories, even if it is not numerical. Credit score is ordinal, because you can rank "good", "fair", and "bad".

Definition 7.8: Nominal

Nominal data is data that has no numeric value and cannot be ordered. Colors such as "red", "blue", "green" cannot be ranked.

Datatypes that are not numeric like ordinal and nominal data are **categorical**.

Some decision trees may or may not require all numerical inputs depending on their implementation. However, in models that use differentiable loss functions (like Linear Regression / Logistic Regression, some forms of decision trees), you need to transform categorical data.

Transforming Ordinal Data: Rank the values (bad = 0, fair = 1, good = 2).

Transforming Nominal Data: Use one-hot encoding.

Definition 7.9: One-Hot Encoding

One-hot encoding is a transformation technique of nominal data. For a feature, it assigns a digit for every possible category that feature can be. For example, for a Color feature, we can assign "red" to the first digit, "blue" to the second digit, and "green" to the third. So, red would be 100, blue would be 010, and green would be 001.

It is important to note that many nominal feature do come in the form of a number, yet they cannot be ordered. An example of this is zip code when estimating house prices. 10018 is a zip code in Manhattan, NY while 98105 is a zip code in Seattle, WA. If we use zip code as a feature of a linear regression model to predict house prices, this does not mean that houses in Seattle are more expensive than those in Manhattan.

7.2.3 Threshold Split for Numeric Features

For numeric features, it doesn't make sense to have individual branches in our tree for every single numerical value that feature could hold. So, we use an inequality instead. One branch will be less than a certain number v , and the other branch can be greater than or equal to. Given a range of numerical values from a single feature, we can choose the best threshold to make this split by calculating the classification error for every possible split that occurs exactly in between each of the possible numerical values v_1, v_2, \dots, v_n and select the threshold yielding the lowest error. For more precision, we can split these inequalities further. There is a limitation to using classification error: Two different splits can give us the same error, so continuous loss functions like entropy loss or Gini impurity loss (which we will not talk about in this text) are used.

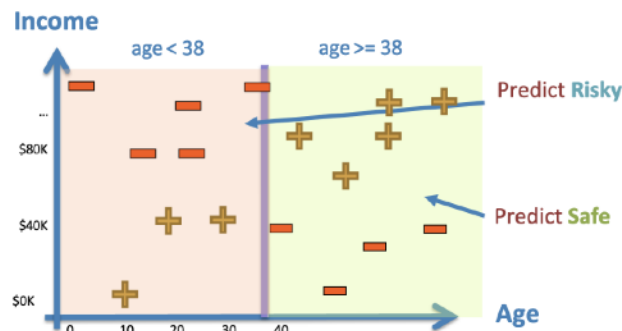


Figure 7.13: For a decision tree, a splitting of a numeric feature will always be along the axis.

7.2.4 Advantages of Decision Trees

1. Decision Trees are easy to interpret
2. They can handle both continuous and categorical variables without preprocessing
3. They do not require normalization
4. They can create non-linear decision boundaries
5. They can handle missing values

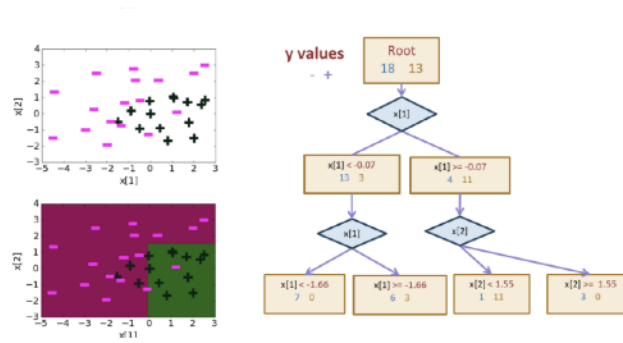


Figure 7.14: If we decide to split further, we continue to split along the axes.

7.2.5 Disadvantages of Decision Trees

1. Deep Decision trees are prone to overfitting. They are not suitable for large datasets for this reason.
2. The decision boundaries are unstable because adding a new datapoint to our trainset can cause the entire tree to regenerate.
3. The decision boundaries must be axis-parallel.

To overcome overfitting, we can implement conditions for early stopping such as establishing a fixed depth length, setting a maximum number of nodes, or halting growth if error does not considerably decrease. We can also **prune** our trees by cutting off some nodes.