

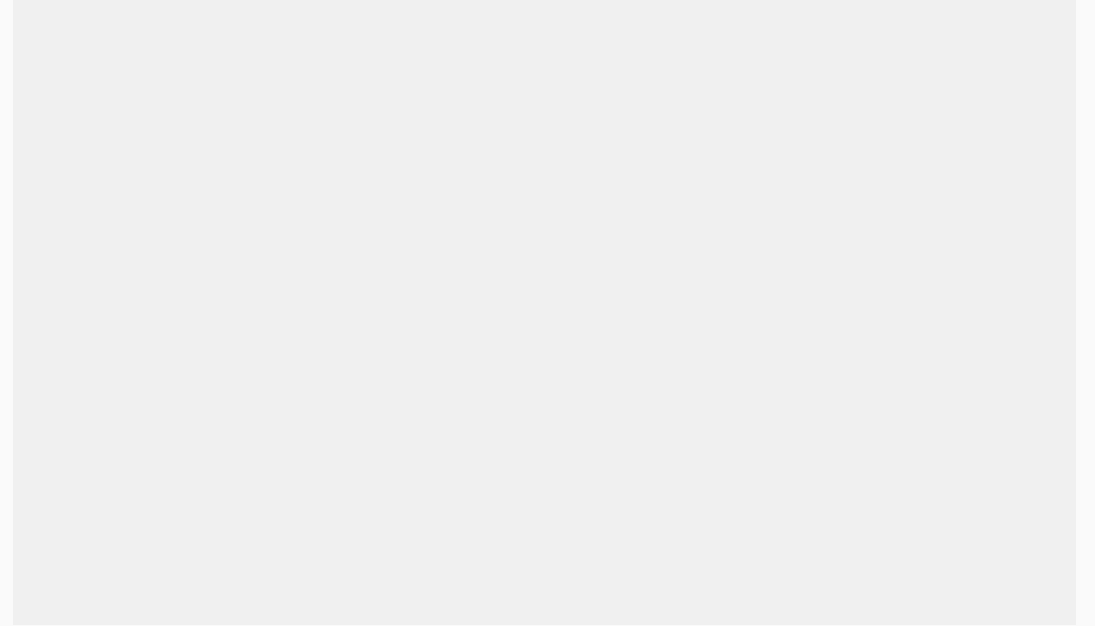
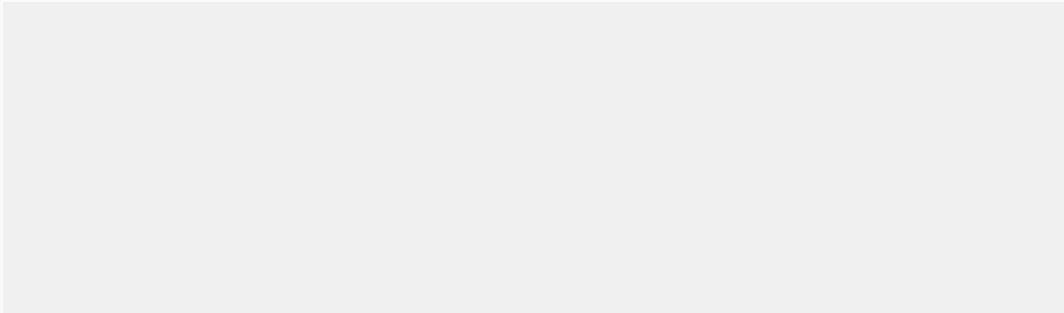
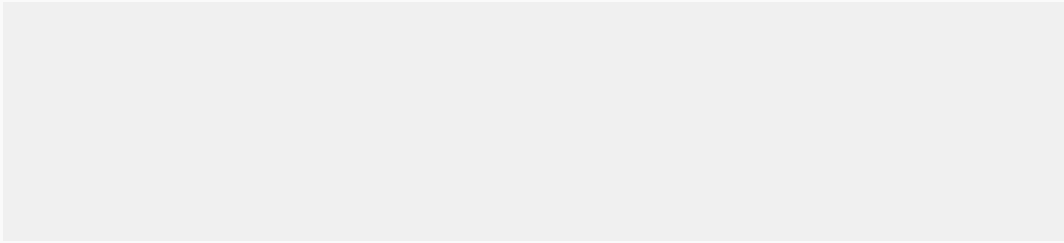
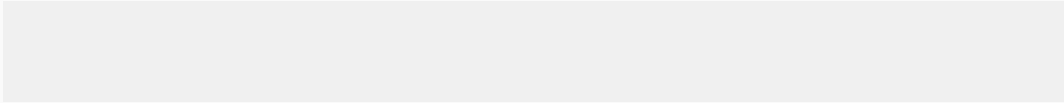
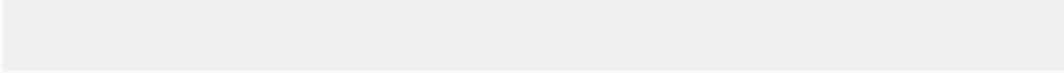
Applied Machine Learning - Classification

Max Kuhn (RStudio)

Outline

- Performance Measures
- OkC Data
- Classification Trees
- More Bagging
- Naive Bayes Models

Load Packages



Measuring Performance in Classification

Illustrative Example



contains another test set example in a data frame called :

```
## Test set example
```


```
## Test set example
```

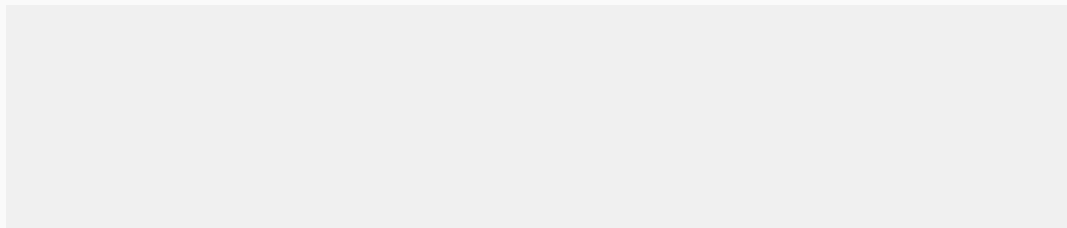
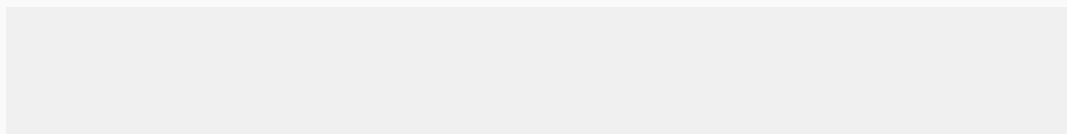
Both `test` and `train` are factors with the same levels. The other two columns represent `prob` and `pred`.

This reflects that most classification models can generate "hard" and "soft" predictions for models.

The class predictions are usually created by thresholding some numeric output of the model (e.g. a class probability) or by choosing the largest value.

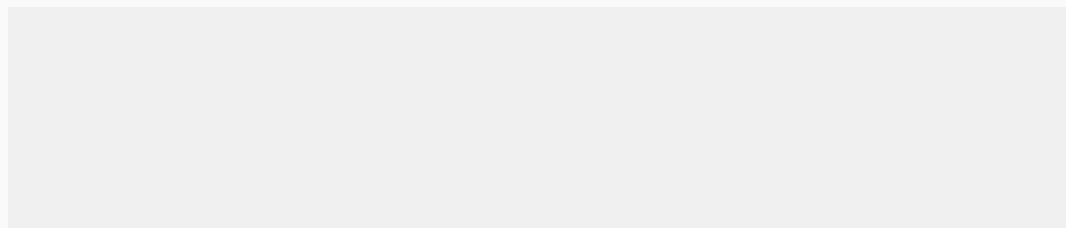
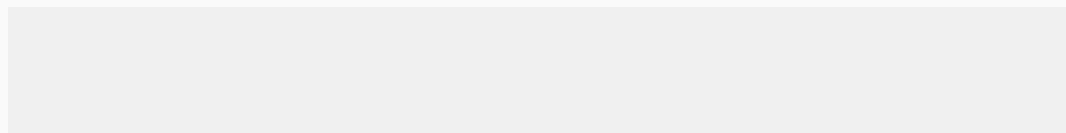
Class Prediction Metrics

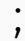
With class predictions, a common summary method is to produce a  which is a simple cross-tabulation between the observed and predicted classes:



These can be visualized using **mosaic plots**.

Accuracy is the most obvious metric for characterizing the performance of models.



However, it suffers when there is a  ; suppose 95% of the data have a specific class. 95% accuracy can be achieved by predicting samples to be the majority class.

Two Classes

There are a number of specialized metrics that can be used when there are two classes. Usually, one of these classes can be considered the `positive` or the `negative`.

One common way to think about performance is to consider false negatives and false positives.

- The sensitivity is the `True Positive Rate` (out of all of the actual positives, how many did you get right?).
- The specificity is the rate of correctly predicted negatives, or `1 - False Positive Rate` (out of all the actual negatives, how many did you get right?).

From this, assuming that `positive` is the event of interest:

$$\text{sensitivity} = 227 / (227 + 31) = 0.88$$

$$\text{specificity} = 192 / (192 + 50) = 0.79$$

Conditional and Unconditional Measures

Sensitivity and specificity can be computed from $\frac{TP}{TP+FN}$ and $\frac{TN}{TN+FP}$, respectively.

It should be noted that these are *conditional* measures since we need to know the true outcome.

The event rate is the $\frac{TP+FN}{TP+FN+TN+FP}$ (or the Bayesian $\frac{P(A)}{P(A)+P(\bar{A})}$). Sensitivity and specificity are analogous to the $\frac{P(B|A)}{P(A)}$ and $\frac{P(\bar{B}|\bar{A})}{P(\bar{A})}$.

There are *unconditional* analogs to the *conditional* measures called the positive predictive values and the negative predicted values.

A variety of other measures are available for two class systems, especially for 2×2 tables.

One thing to consider: what happens if our $\frac{TP}{TP+FN}$ is low? $\frac{TN}{TN+FP}$ is low?

Changing the Probability Threshold

For two classes, the 50% cutoff is customary; if the probability of class #1 is $\geq 50\%$, they would be labelled as class \#1 .

What happens when you change the cutoff?

- Increasing it makes it harder to be called class \#1
fewer predicted events, specificity \uparrow , sensitivity \downarrow
- Decreasing the cutoff makes it easier to be called class \#1
more predicted events, specificity \downarrow , sensitivity \uparrow

With two classes, the

ROC curve

can be used to estimate performance using a combination of sensitivity and specificity.

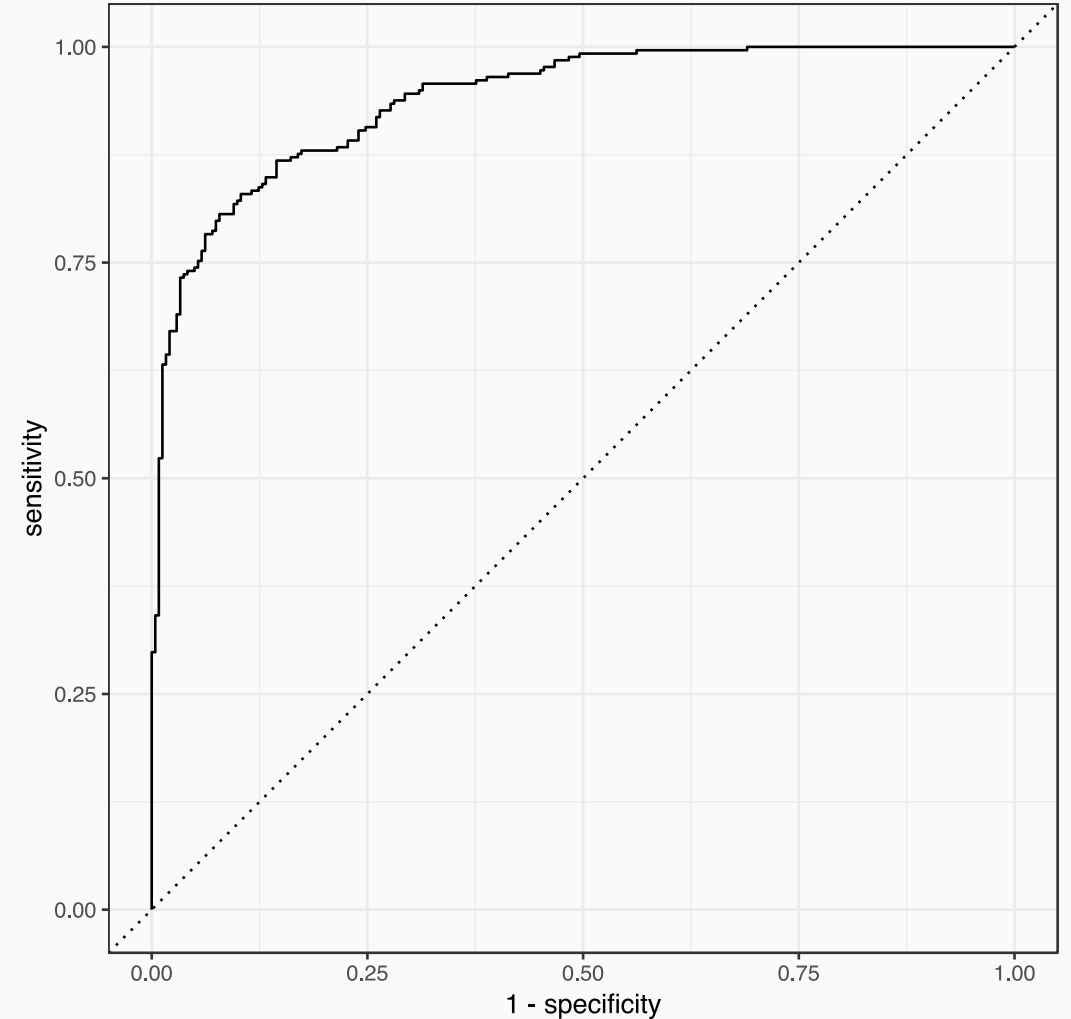
To create the curve, many alternative cutoffs are evaluated.

For each cutoff, we calculate the sensitivity and specificity.

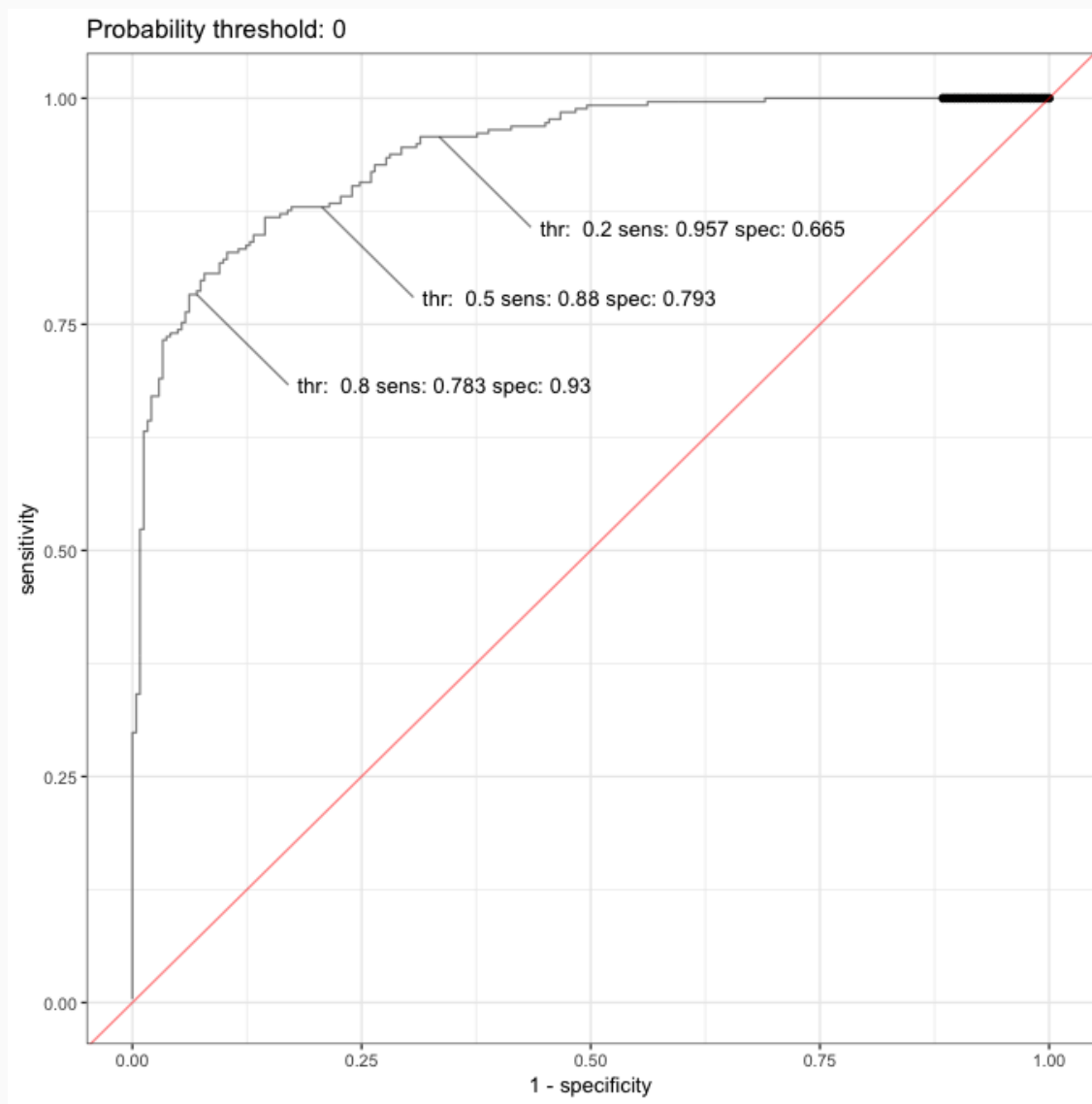
The ROC curve plots the sensitivity (eg. true positive rate) versus $1 - \text{specificity}$ (eg. the false positive rate).

The area under the ROC curve is a common metric of performance.

The Receiver Operating Characteristic (ROC) Curve



Changing the Threshold



The Receiver Operating Characteristic (ROC) Curve

The ROC curve has some major advantages:

- It can allow models to be optimized for performance before a definitive cutoff is determined.
- It is robust to class imbalances; no matter the event rate, it does a good job at characterizing model performance.
- The ROC curve can be used to pick an optimal cutoff based on the trade-offs between the types of errors that can occur.

When there are two classes, it is advisable to focus on the area under the ROC curve instead of sensitivity and specificity.

Once an acceptable model is determined, a proper cutoff can be determined.

OkC Data

These data contains several types of fields:

- a number of open text essays related to interests and personal descriptions
- single choice type fields, such as profession, diet, gender, body type, etc.
- multiple choice data, including languages spoken, etc.

We will try to predict whether someone has a profession in the STEM fields (science, technology, engineering, and math) using a random sample of the overall dataset.

The data are included in the workshop's GitHub repo:



Dealing with Class Imbalances

In our data, only 18.2% of the profiles are in STEM professions. This complicates the analysis since many models will overfit to the majority class.

There are two main strategies to deal with this:

- **Cost-sensitive models** where a higher cost is attached to the minority classes. In this way, the fitting process puts more emphasis on those samples.
- **Resampling** that modify the rows of the data to re-balance the training set.

Cost-sensitive models tend to only produce hard classifications so we will focus on the latter.

Class Imbalance Sampling

There are a variety of methods for subsampling the data. Some exclude or replicate rows in the training set while others try to add new data points to balance the classes.

The simplest method for dealing with the problem is to randomly sample the data to make the number of STEM and non-STEM profiles the same.

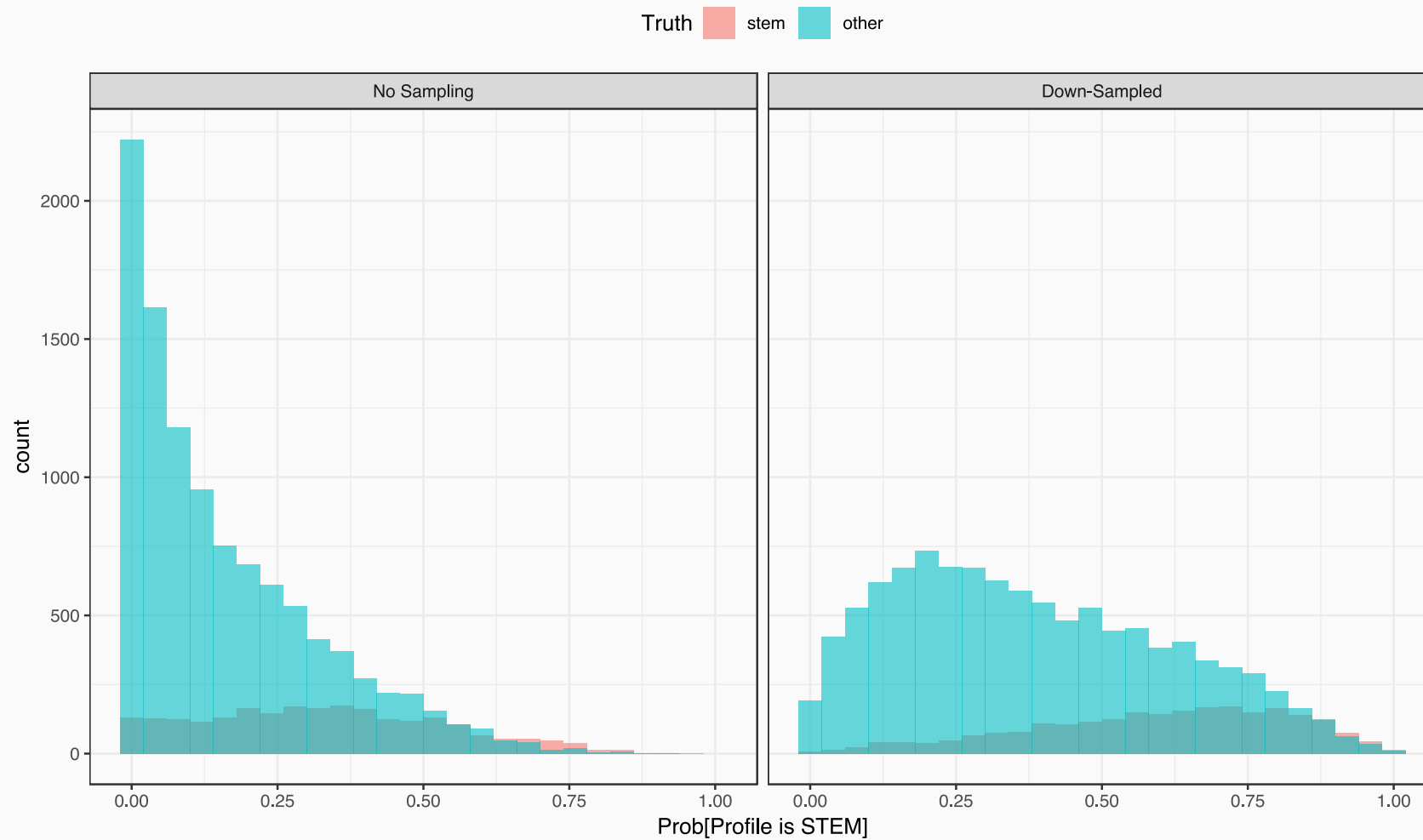
While it seems like throwing away most of the data is a bad idea, it tends to produce less pathological distributions of the class probabilities and can help to improve the ROC curve.

It is critical that:

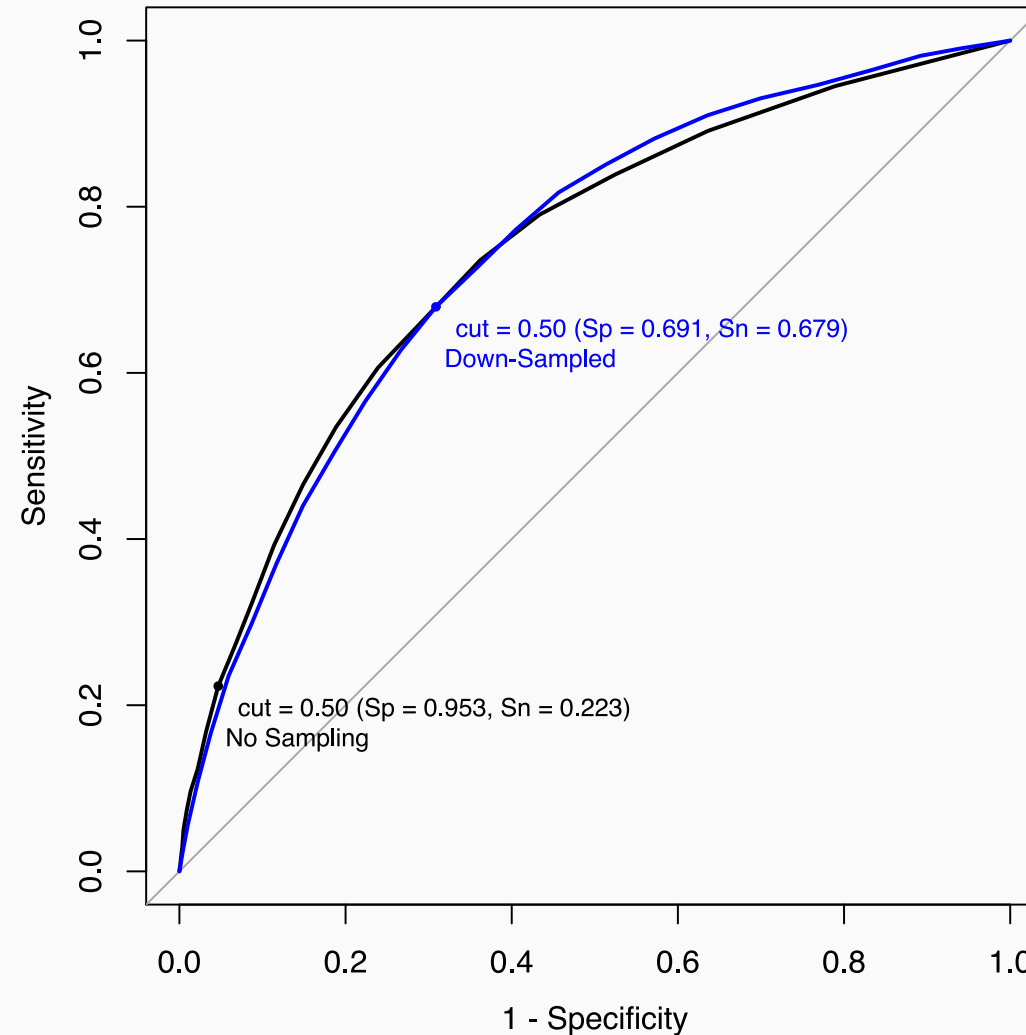
- the sampling should be done on the training set. Otherwise, the performance estimates can be optimistic.
- these sampling methods take place on the analysis set and not the assessment set.

Note that for a simple logistic regression model, this mainly has the effect of changing the intercept.

Calibration Effect (Test Set Example)



Calibration Effect (Test Set Example)



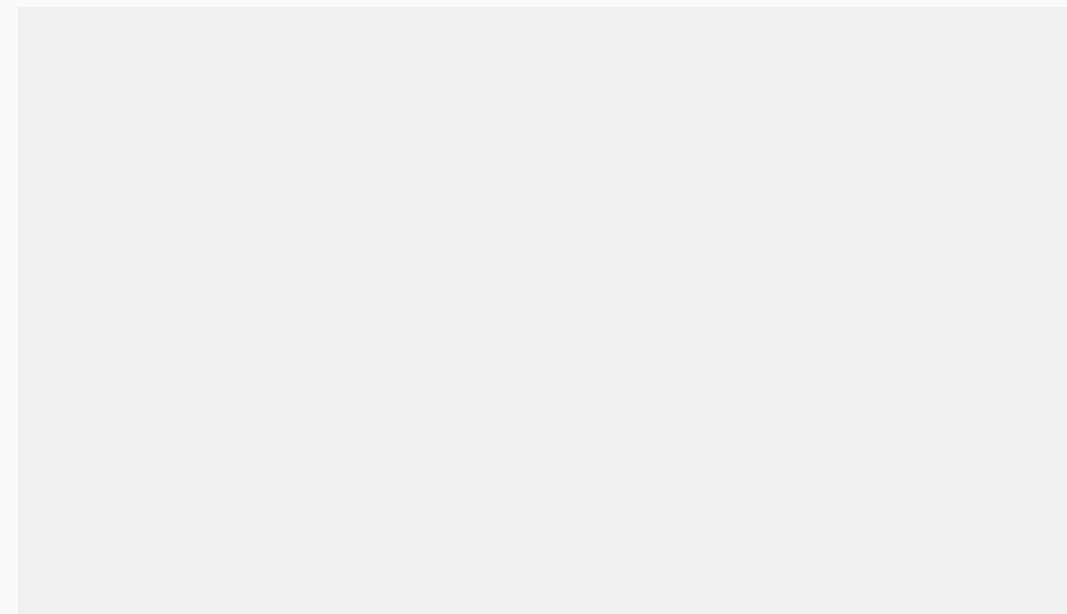
Resampling and Analysis Strategy

If we down-sample the data, the analysis set will consist of 1458 profiles (equally balanced). We'll again use 10-fold CV to resample the data.

Within each resample, the analysis data are down-sampled and the assessment sets are left alone.

The number of STEM profiles held-out would be about 72 and this should be sufficient to compute sensitivity.

The models will be optimized on the area under the ROC curve.



Classification Trees

Classification Trees

Tree-based classifiers conduct searches of the predictors to find the best split of the data to create two subsets.

"Best", in most cases, means that the class distribution is as as possible in the subsets (i.e. is mostly one class).

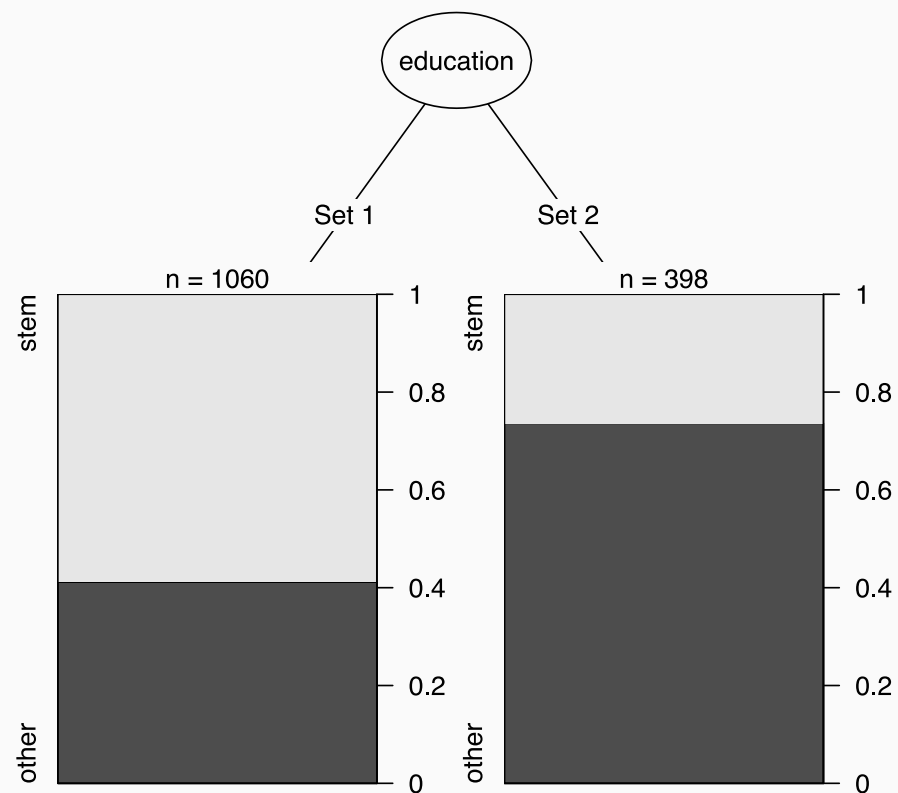
There are many different types of classification trees, including conditional inference trees, C5.0, Bayesian additive regression trees, globally optimal trees, CART, and others.

We'll focus on CART via the package for these notes.

Two Possible Splits - Which One is Better?

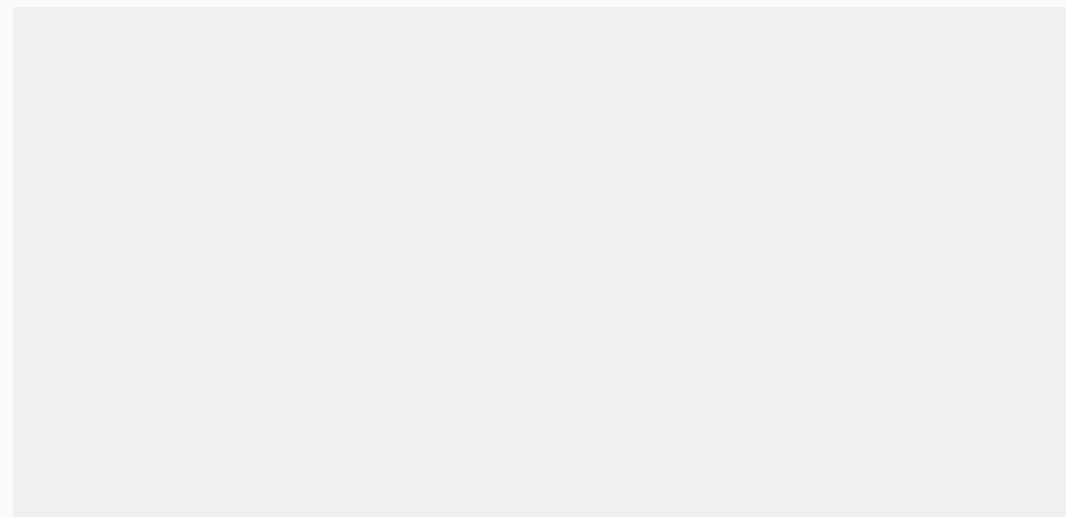
The data have been down-sampled so that classes have equal frequencies.

A Better Split

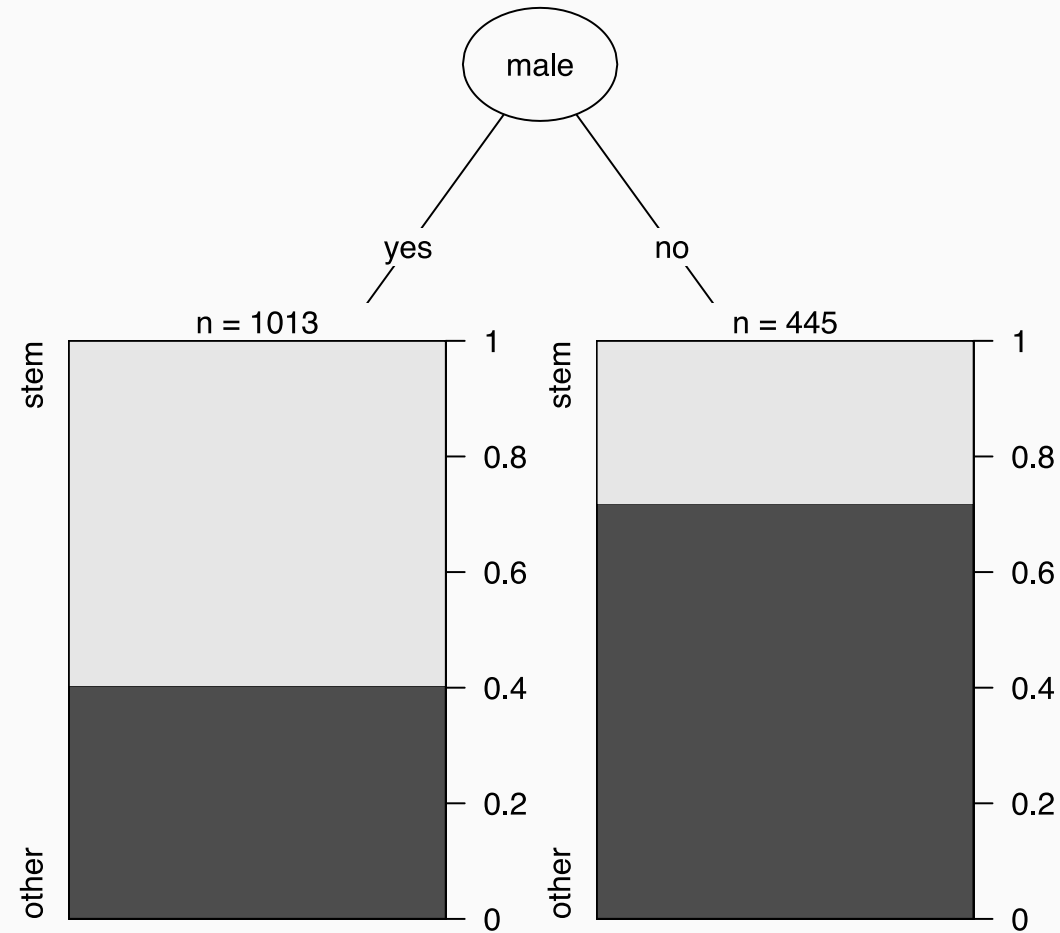


The two sets were derived by the model and are not listed here due to their sizes.

"Set 1" includes



The Official First Split

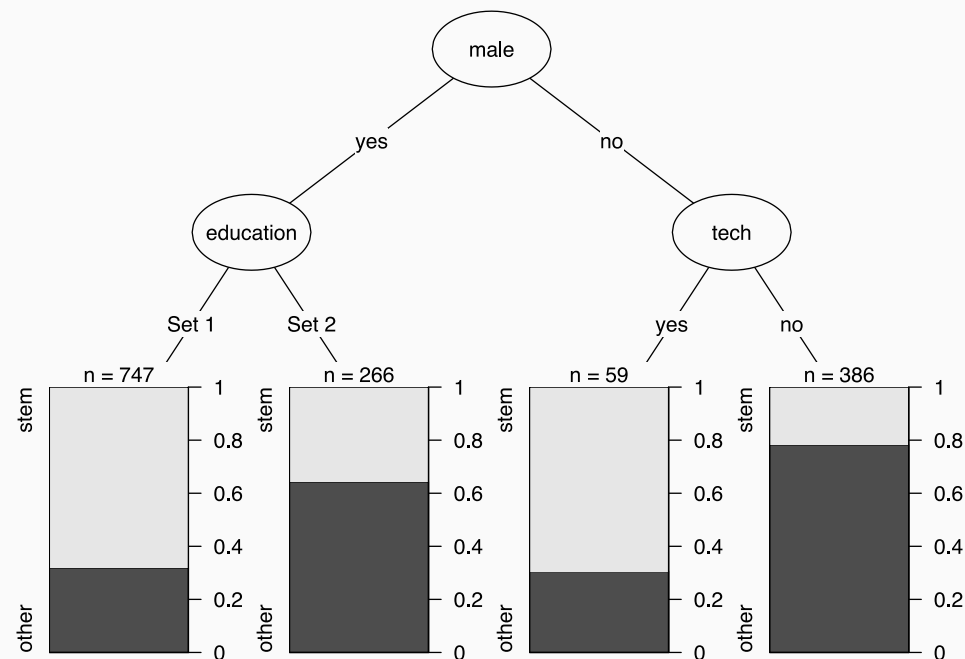


The Next Two splits

Once the initial split is made, the model will then split the resulting two data sets using new searches in those .

The process continues until there are not enough data points left to accurately split or a pre-defined split limit has been reached.

This is the process and, once complete, most tree-based models begin to the trees using some method that balances complexity with performance.

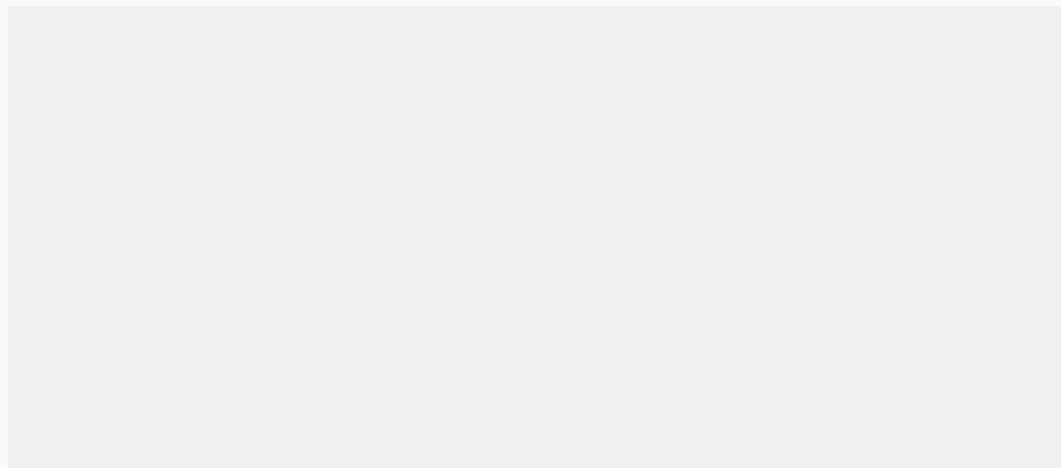


Classification Trees

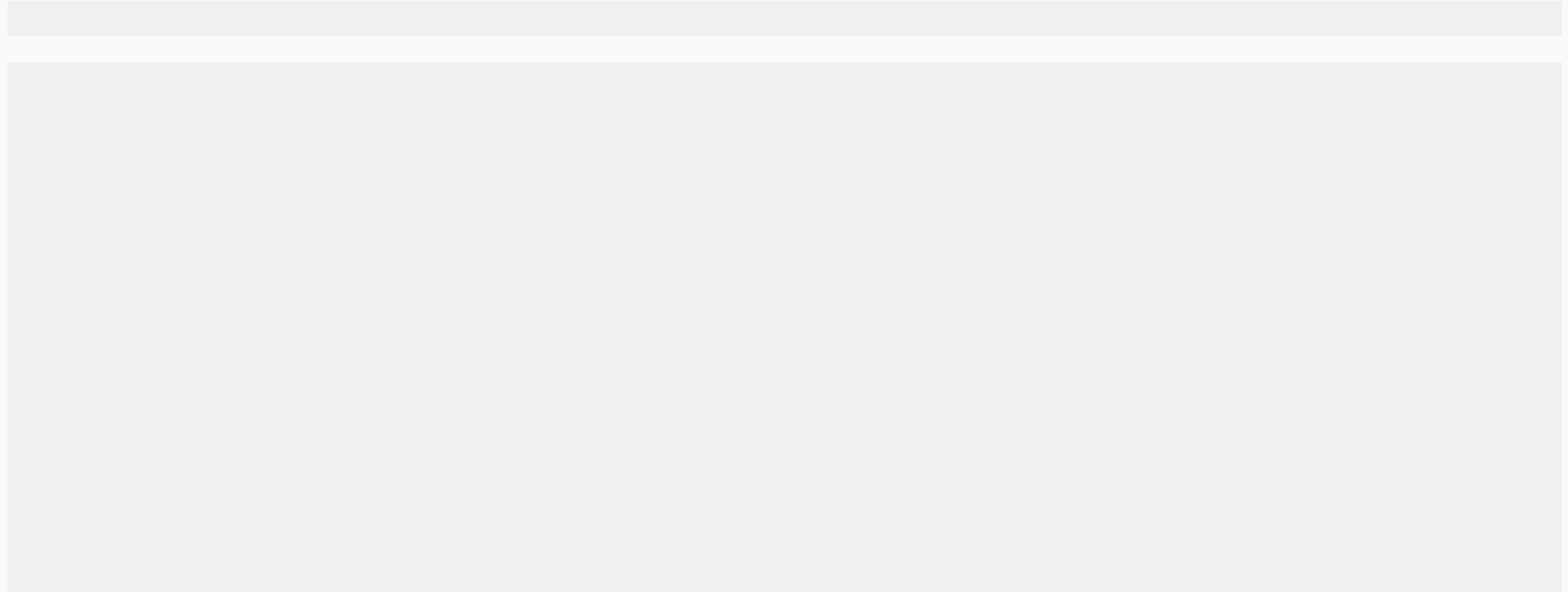
`gbm` contains multiple method for training CART models.

We'll go with `gbm` which uses `cv.gbm` as the tuning parameter.

Also, we'll use the non-formula interface so that the factor predictors are not converted to indicator variables.

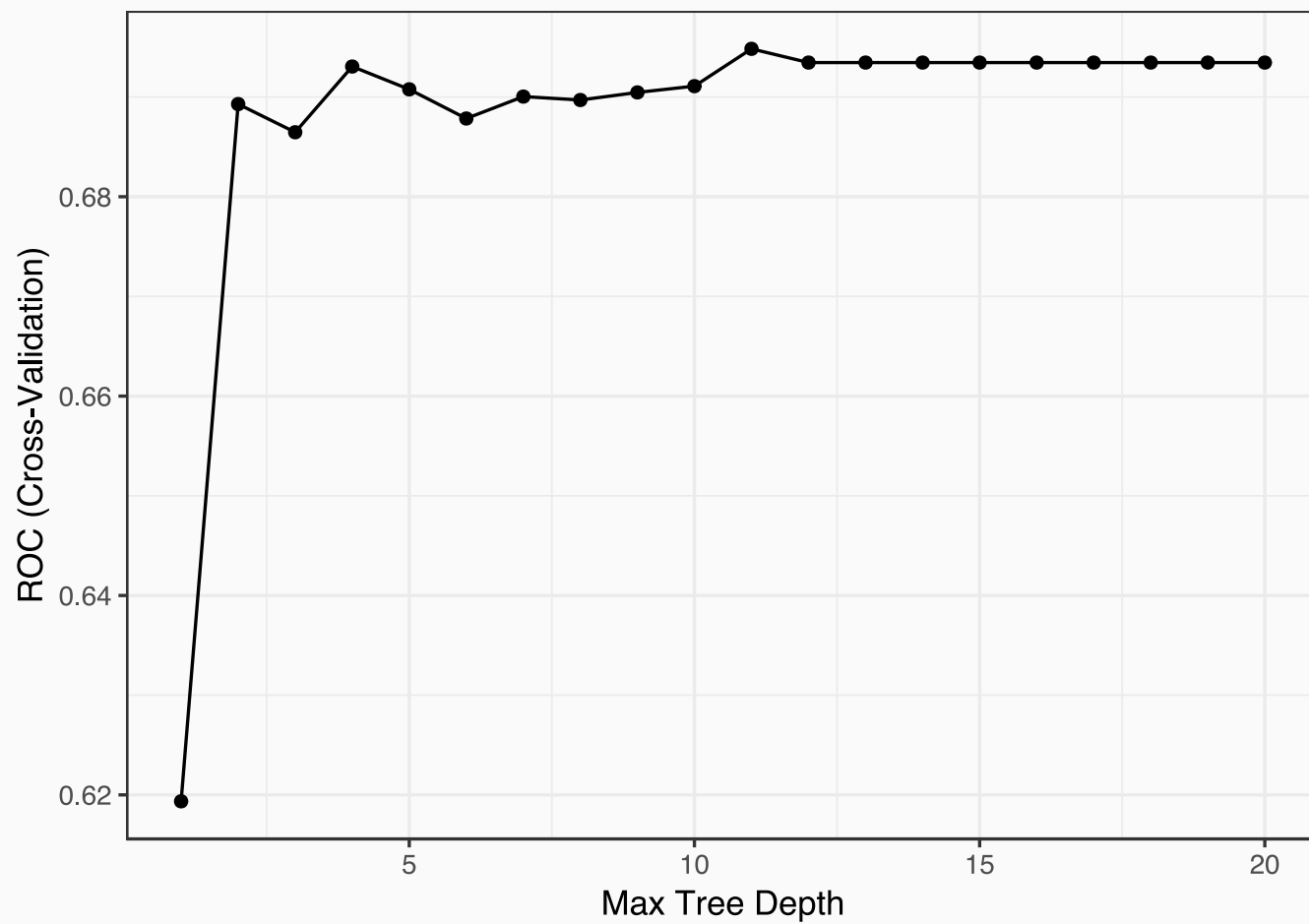


CART Model



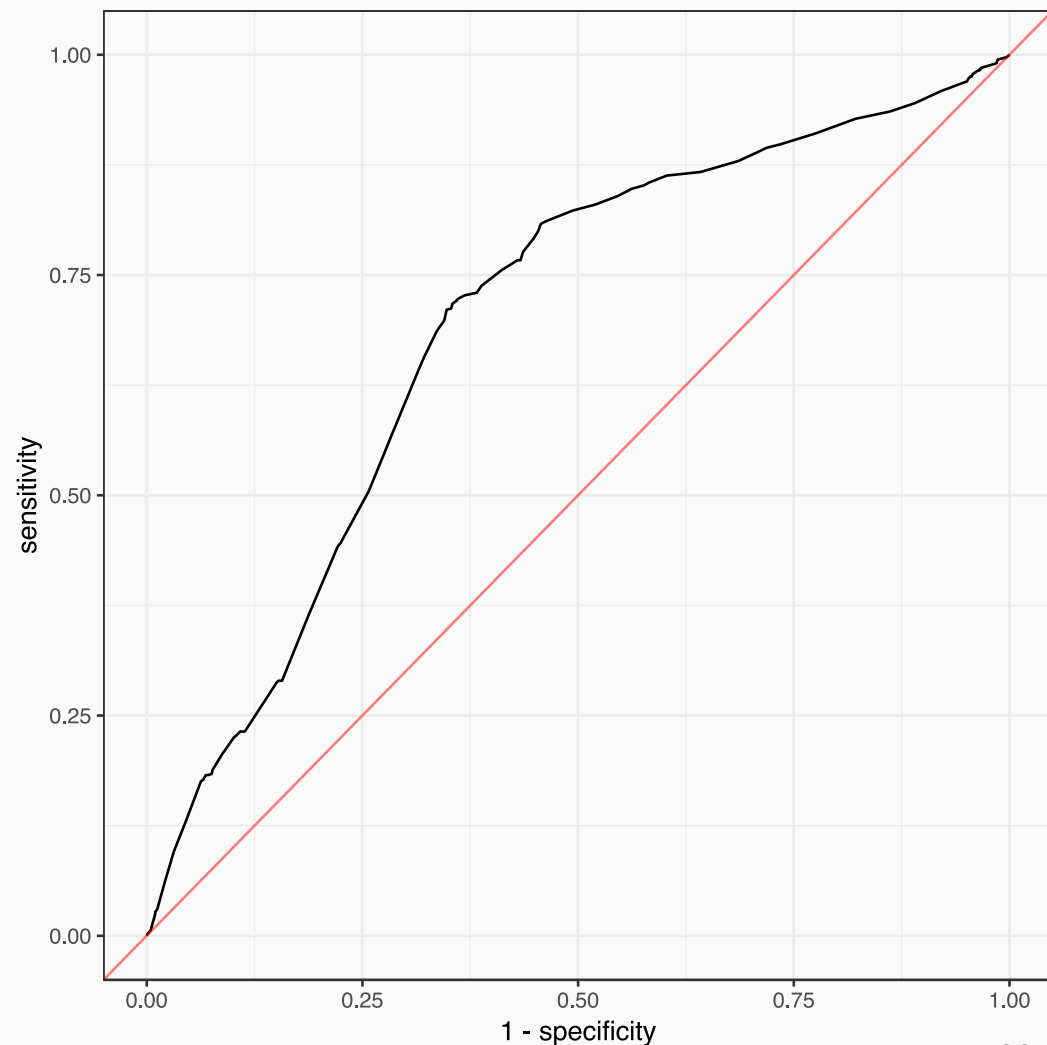
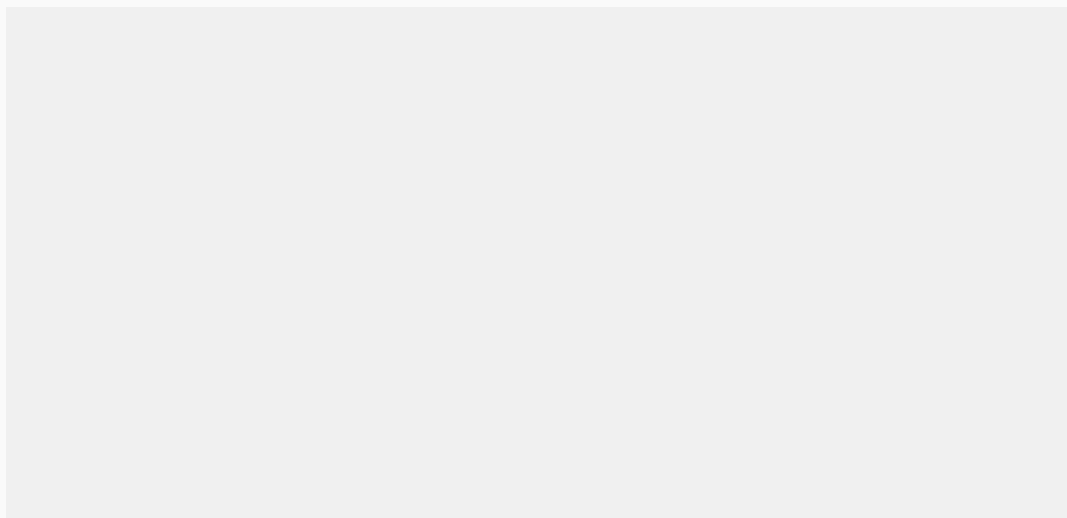
Note that there are 10 terminal nodes and thus 10 possible probability values.

CART Resampling Profile



Classification Tree Average ROC Curve

ROC curves will be created for each model so a convenience function will be used repeatedly to create an ROC curve:

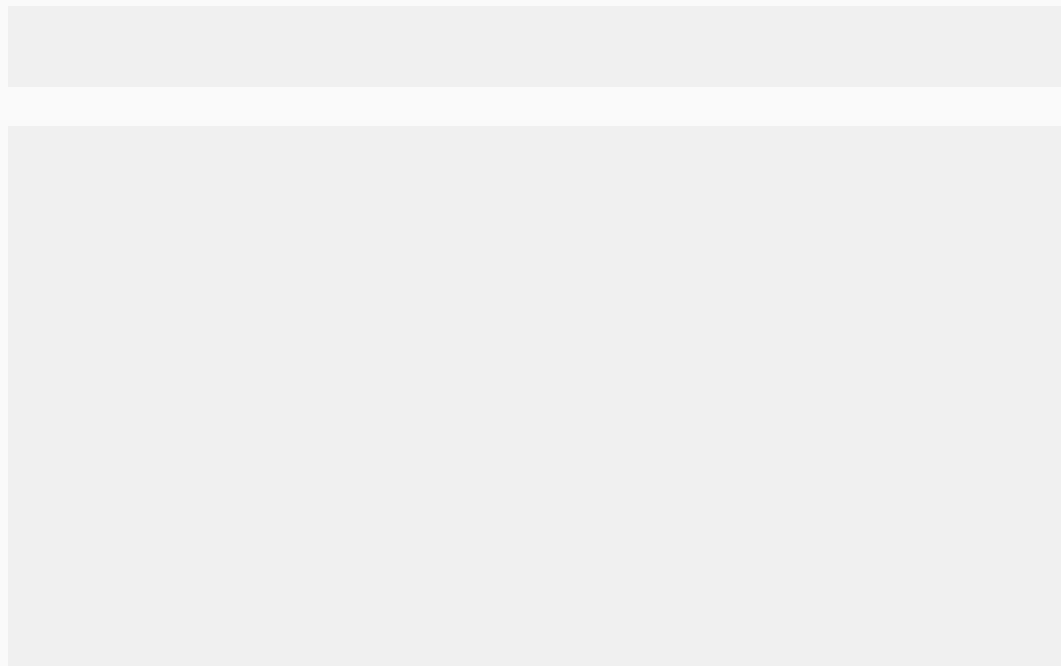


Resampled Confusion Matrix

Since there are 10 different assessment sets, separate confusion matrices can be constructed for each.

`confusionMatrix` has its own `aggregate` method that can show aggregations of these matrices.

By default, the overall rates are shown in each cell.



Variable Importance Scores

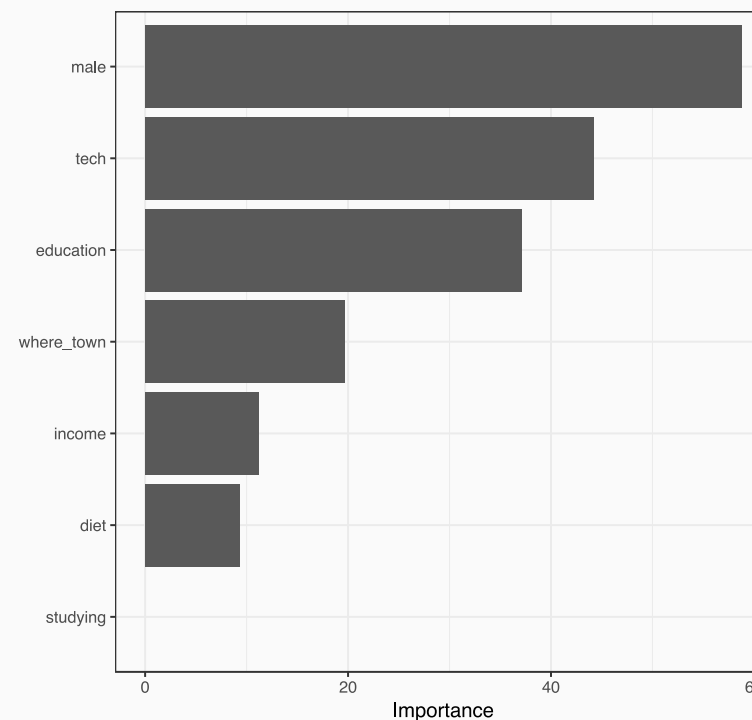


CART tracks the decrease in impurity in the nodes when splits are made.

These can be aggregated over the tree to produce another general importance score.

can also keep track of splits not used in the model (e.g. surrogate, completing) that can be used when there are missing predictor values. The method has options to turn these on/off.

Turning these off gives only the splits uses in the official tree.



Hands-On: I Don't Want to be a Dummy!

Take the previous code and `get_dummies()` to create the model. For `get_dummies()`, the formula method creates dummy variables for predictors that are factors.

Is there any difference in performance?

Is the final model affected? How?

Take 15 min to answer these questions.

Bagging (Again)

Instability of Trees

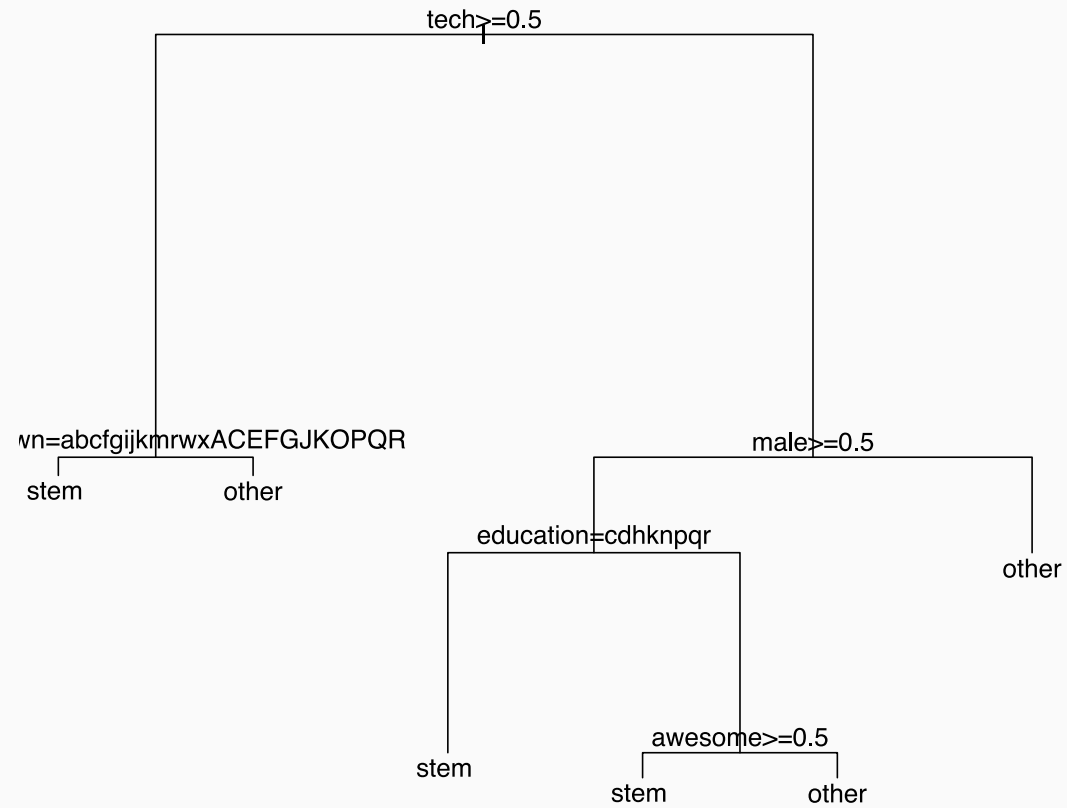
Many types of trees are `unstable` in that they have high variance; if the data are slightly changed, a large impact can be seen on the structure of the model.

This is generally bad and might make you question the interpretability of trees.

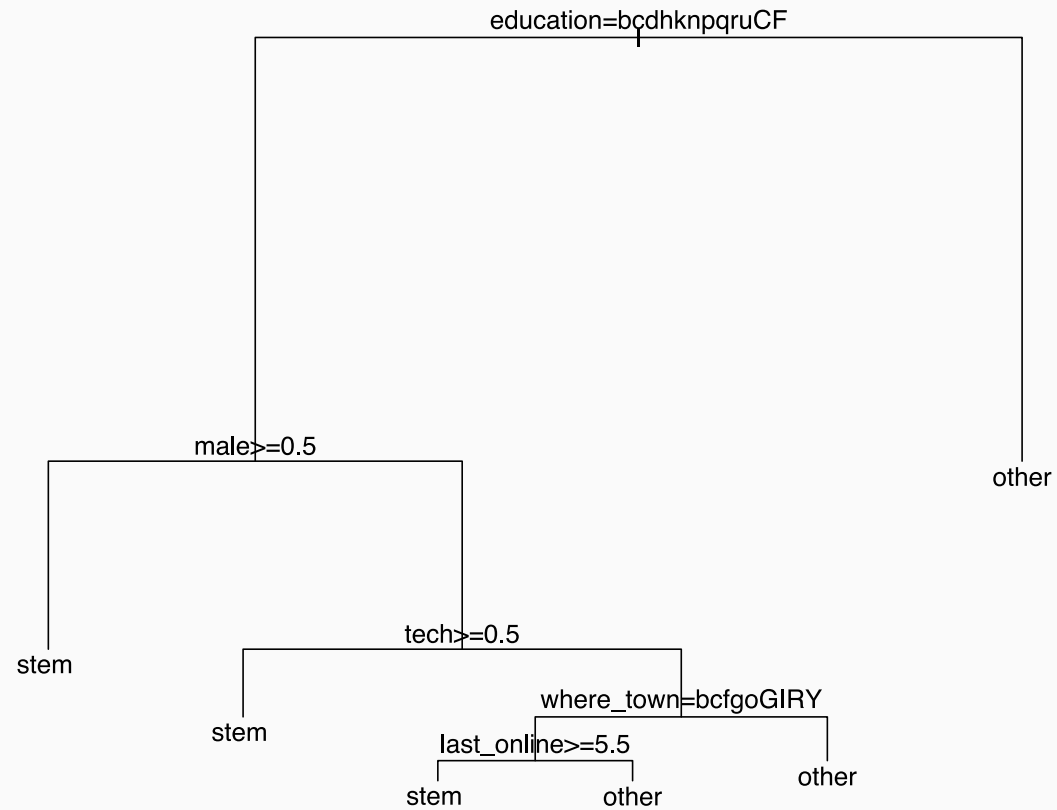
For example, what happens if we were to build our model on bootstrap samples of the training set?

In the three following plots, the maximum tree depth was capped at 5 to make the trees easier to visualize.

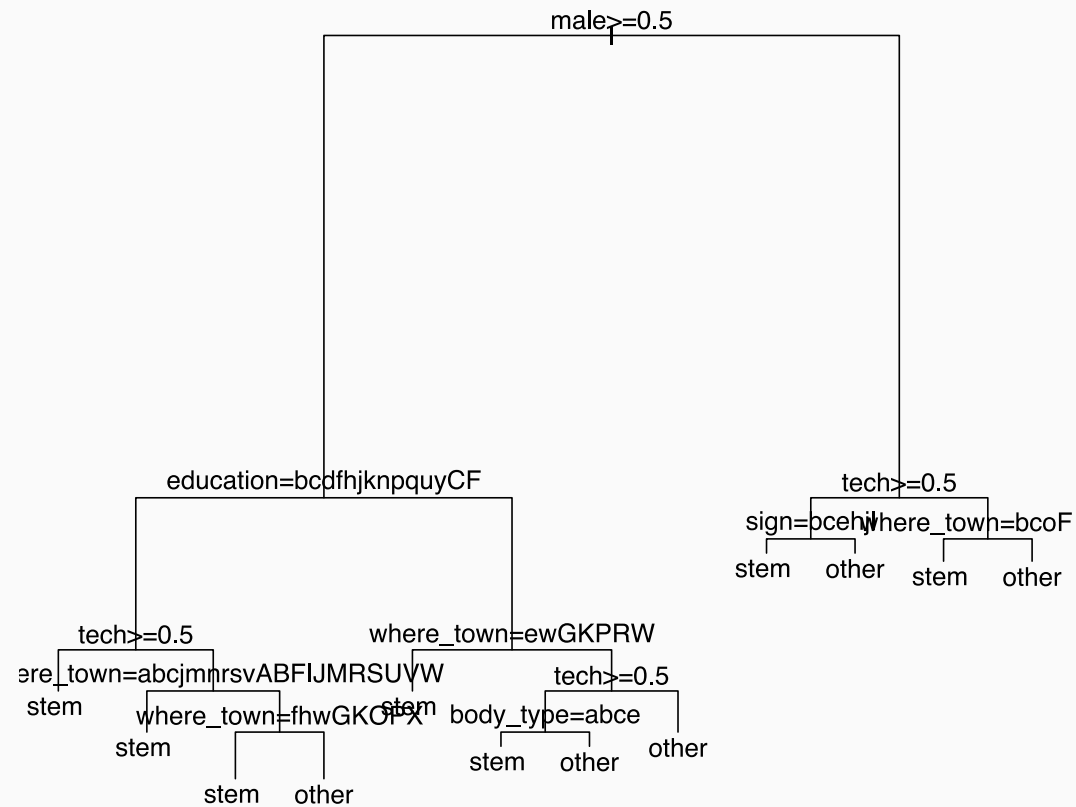
Bootstrap Sample #1



Bootstrap Sample #2



Bootstrap Sample #3



Lemonade from Lemons

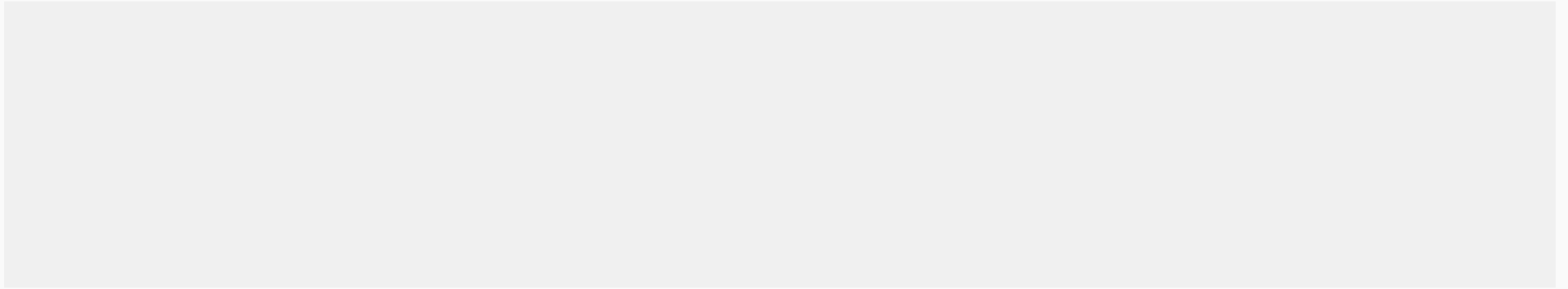
The thing about this instability is that it makes trees great candidates for ensemble methods.

Since ensembles use multiple models, they are only effective when the constituent models are ; otherwise the same predictions are averaged.

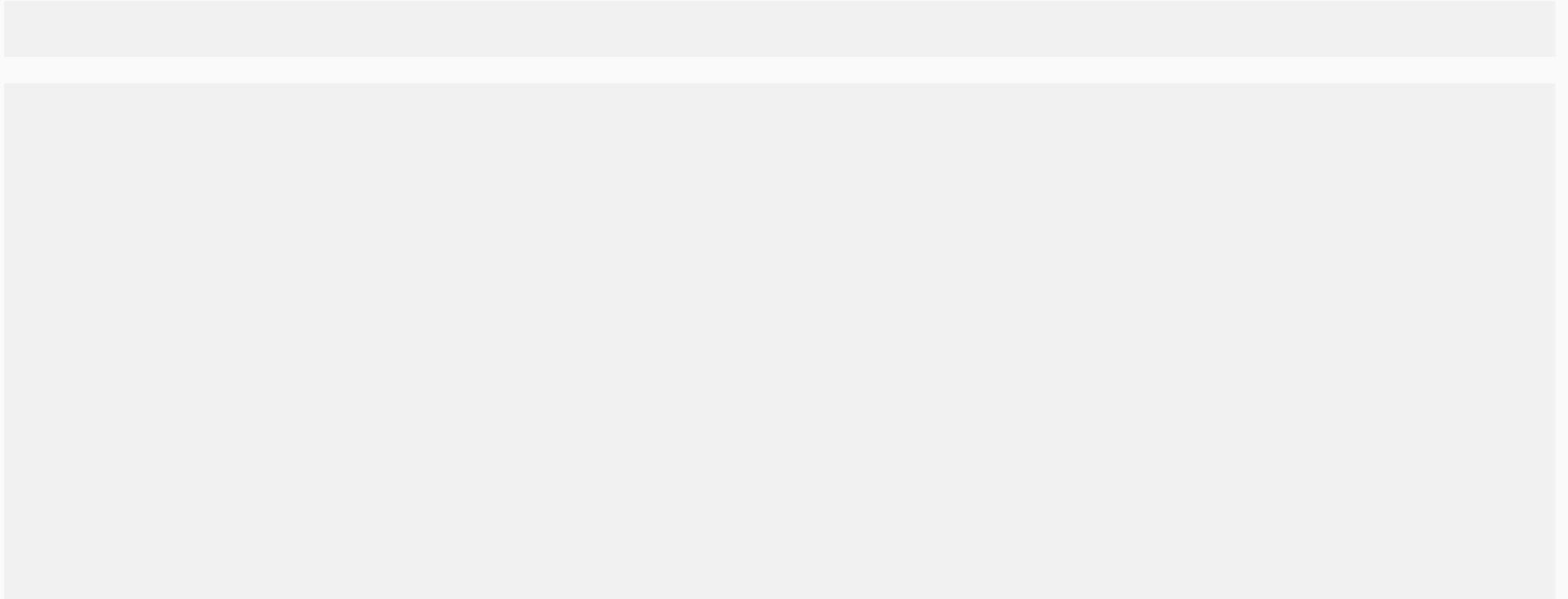
Let's bag the CART trees by using bootstrap samples of the training set and growing the largest possible tree.

wraps using . We will use the default ensemble size of 25 models.

Bagging CART Trees

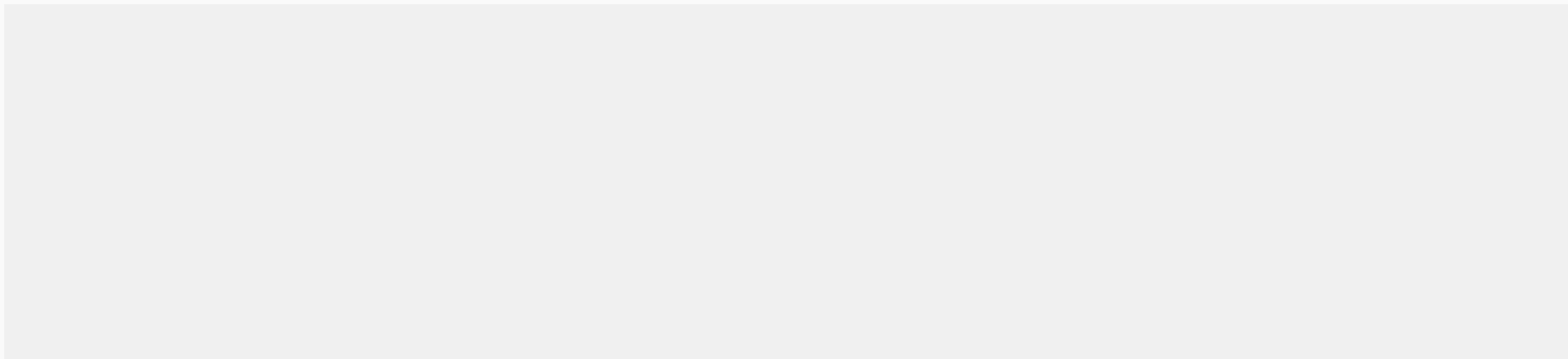
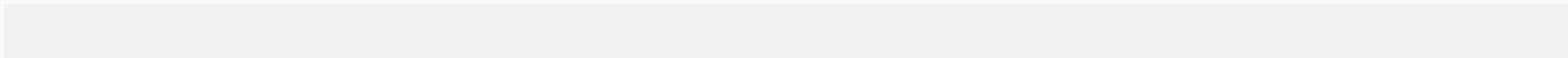


Bagged CART Results

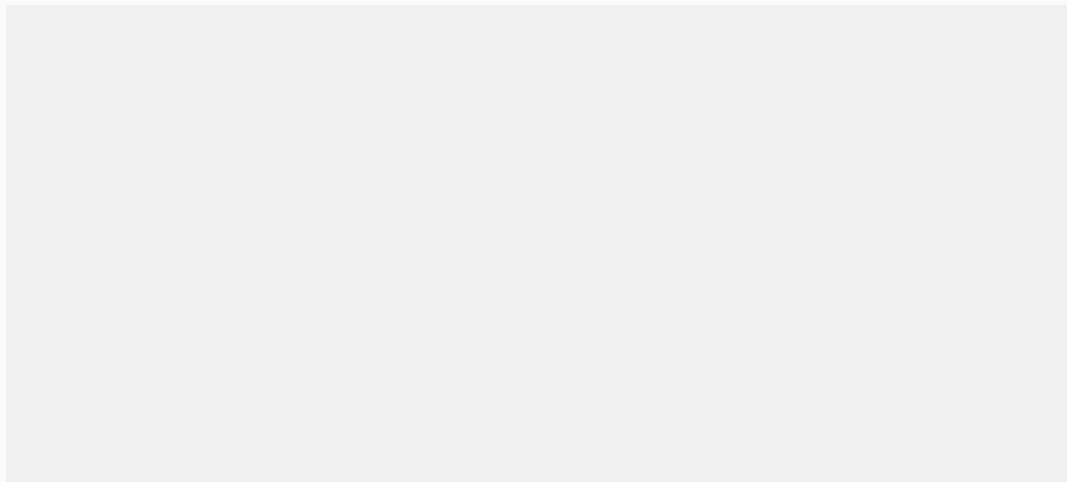


Resampled Confusion Matrix

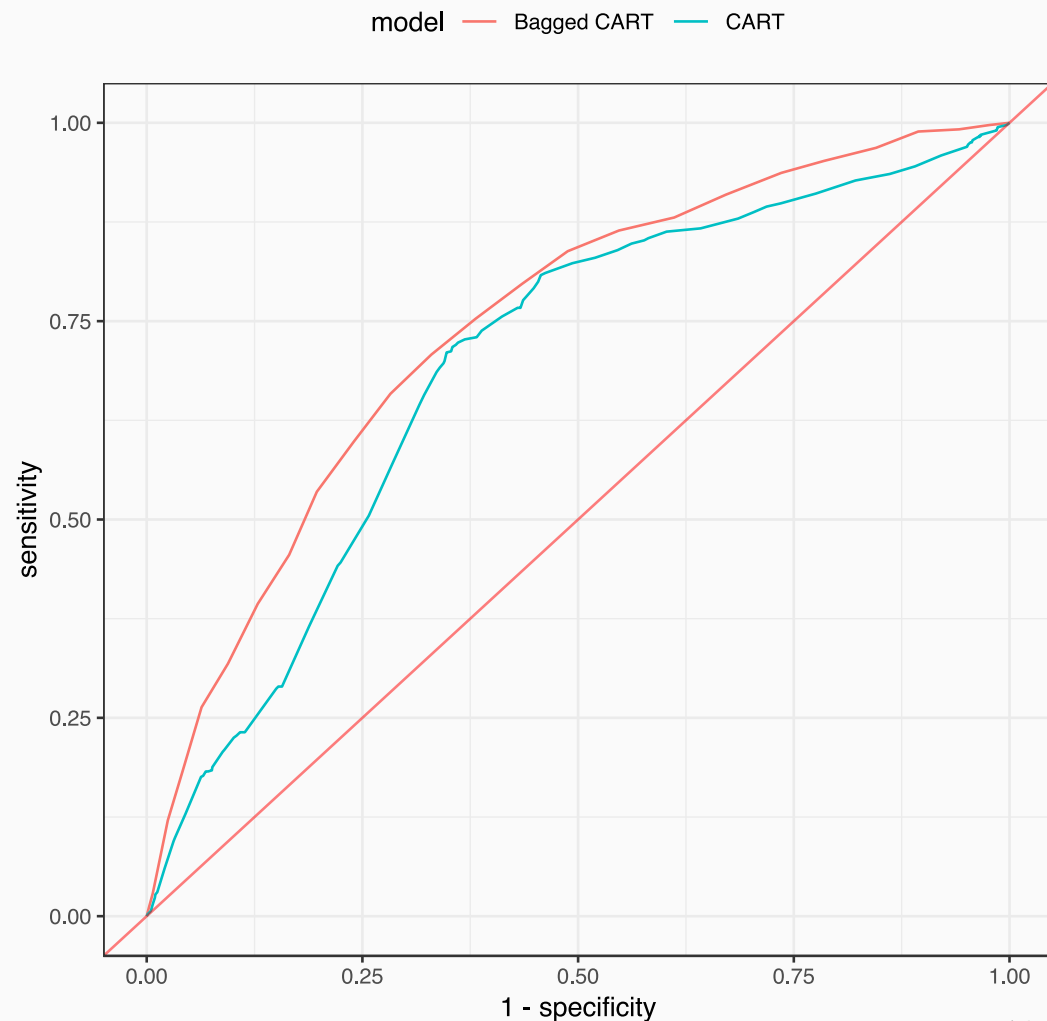
As measured by the default cutoffs, there is some increase in accuracy that is achieved by improving the specificity (27% versus CART's 30.1%).



Bagged Classification Tree Average ROC Curve



Although the bagged model's curve is uniformly better than the single tree, their performance is very similar for the cutoffs closest to upper left-hand corner.

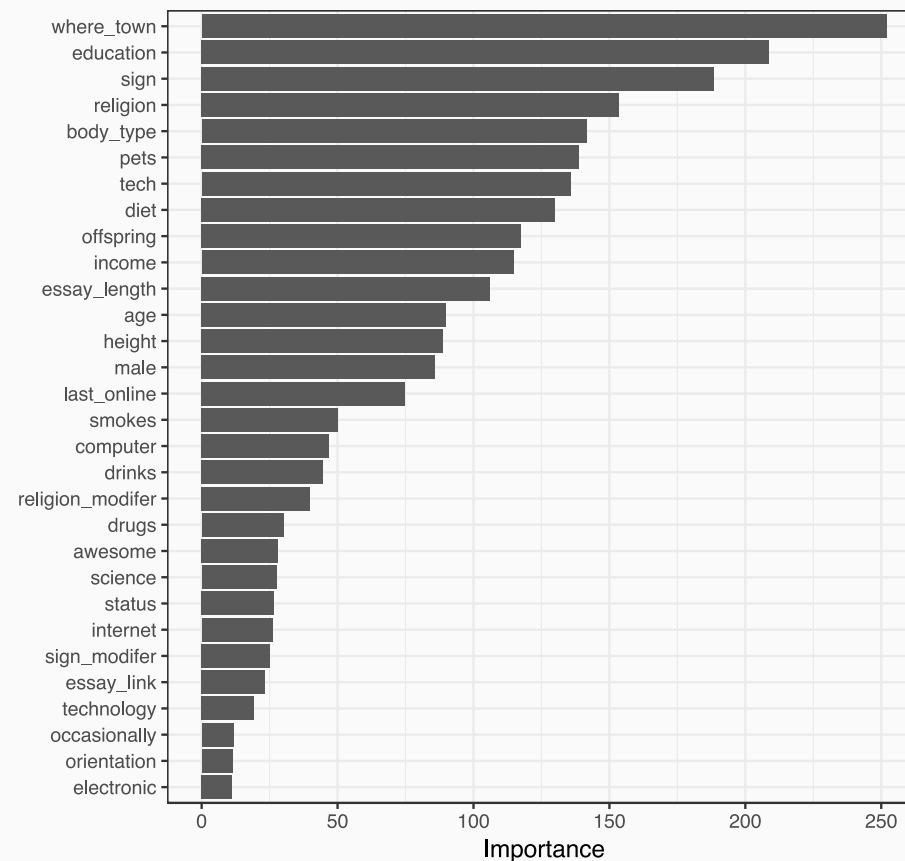


Variable Importance Scores



When bagging, the CART importances scores for each tree are aggregated across trees.

Many more predictors are used in this model.



Hands-On: How Many Trees?

As previously mentioned, `rf` uses `max_depth` to create the model.

Look at the help function to determine which `n_estimators` argument controls the number of bootstraps.

How can we make `rf` use a different value? (hint: `n_estimators`)

Does changing this affect the area under the ROC curve?

Take another 10 mins.

Other Ensemble Methods

There are a variety of other methods for creating ensembles

- `randomForest` are just like bagging but the trees are made more diverse by randomly sampling a subset of predictors to be used in each split. (`randomForest`)
- `boost` fits a sequence of trees and modifies the case-weights of each data point to increase diversity. (`boost`, `adaboost`)
- `randomForest` is PCA signal extraction on a random subset of the data prior to creating the trees. (`randomForest`)
- Regression `boost` adjust the outcome data over a sequence of models. (`boost`)
- `merf` is an ensemble method where different types of models can be blended together through averaging. (`merf`)

R has multiple implementations of these methods.

Bayes' Rule

Naive Bayes Models

This classification model is motivated directly from statistical theory based on Bayes' Rule:

In English:

Given our predictor data, what is the probability of each class?

The _____ is the prevalence that was mentioned earlier (e.g. the rate of STEM profiles). This can be estimated or set.

Most of the action is in _____, which is based on the observed training set.

Predictions are based on a blend of the training data and our _____ about the outcome...

So Why is it Naive?

Determining `probabilities` can be very difficult without strong assumptions because it measures the `joint probability` of all of the predictors.

- For example, what is the correlation between a person's essay length and their religion?

To resolve this, `Naive Bayes` assumes that all of the predictors are `independent` and that their probabilities can be estimated separately.

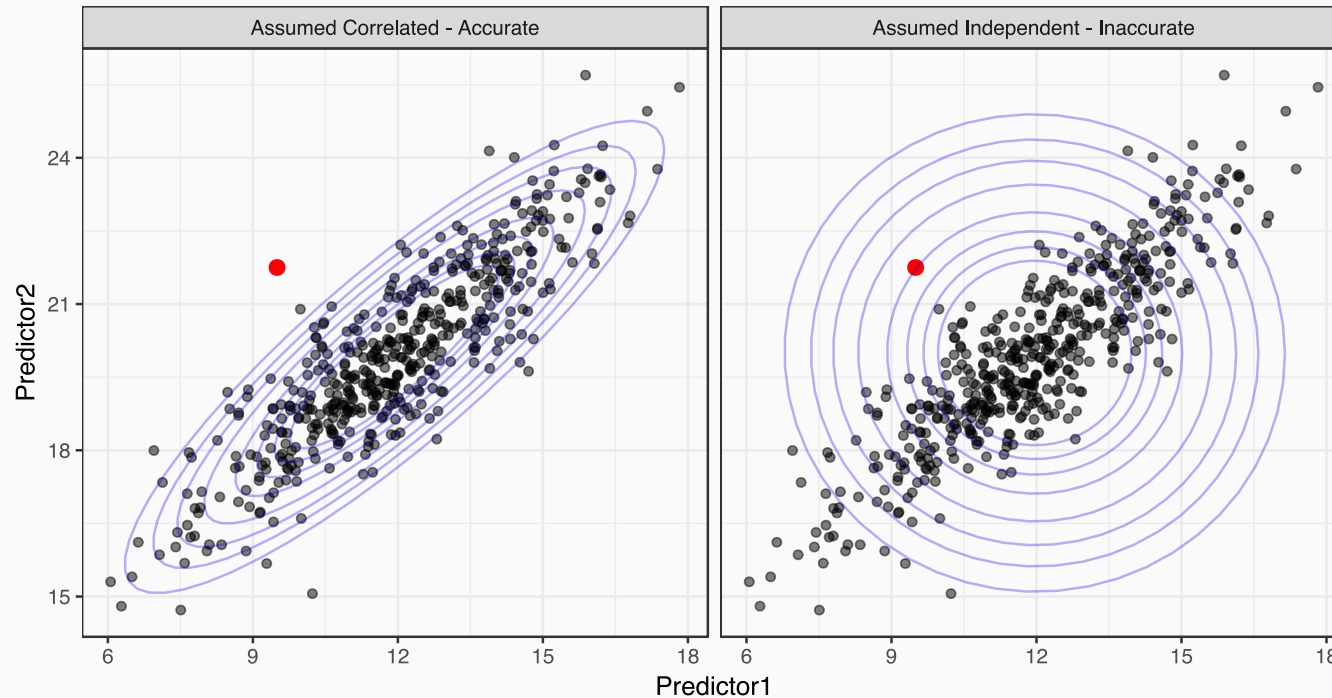
The joint probability is then the product of all of the individual probabilities (an example follows soon).

This assumption is almost certainly bogus but the model tends to do well despite this.

The Effect of Independence

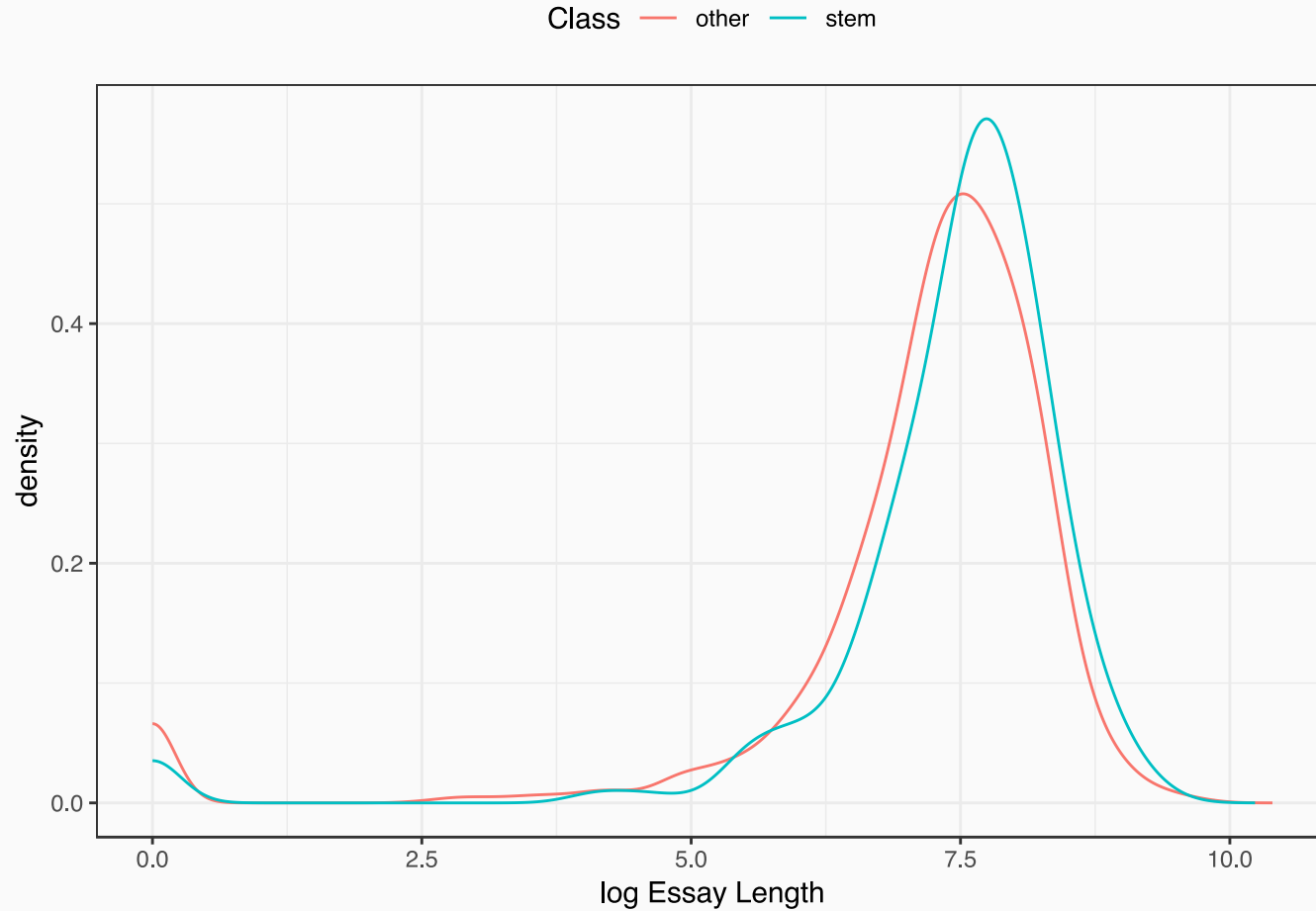
The probability contours assume multivariate normality with different assumptions.

Suppose the red dot is a new sample.

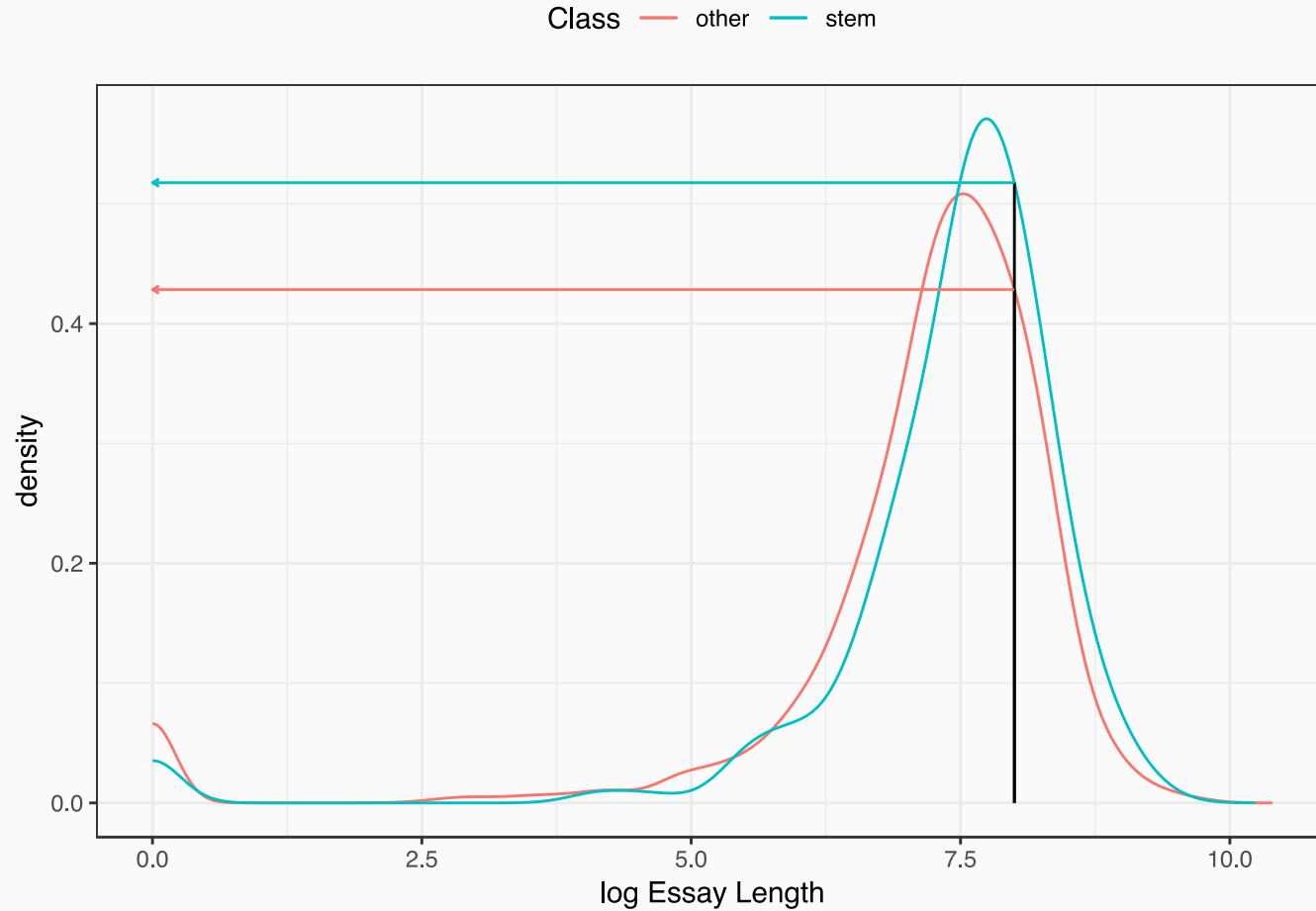


Probability of the red point: 0.0000066 (accurate) and 0.013 (inaccurate).

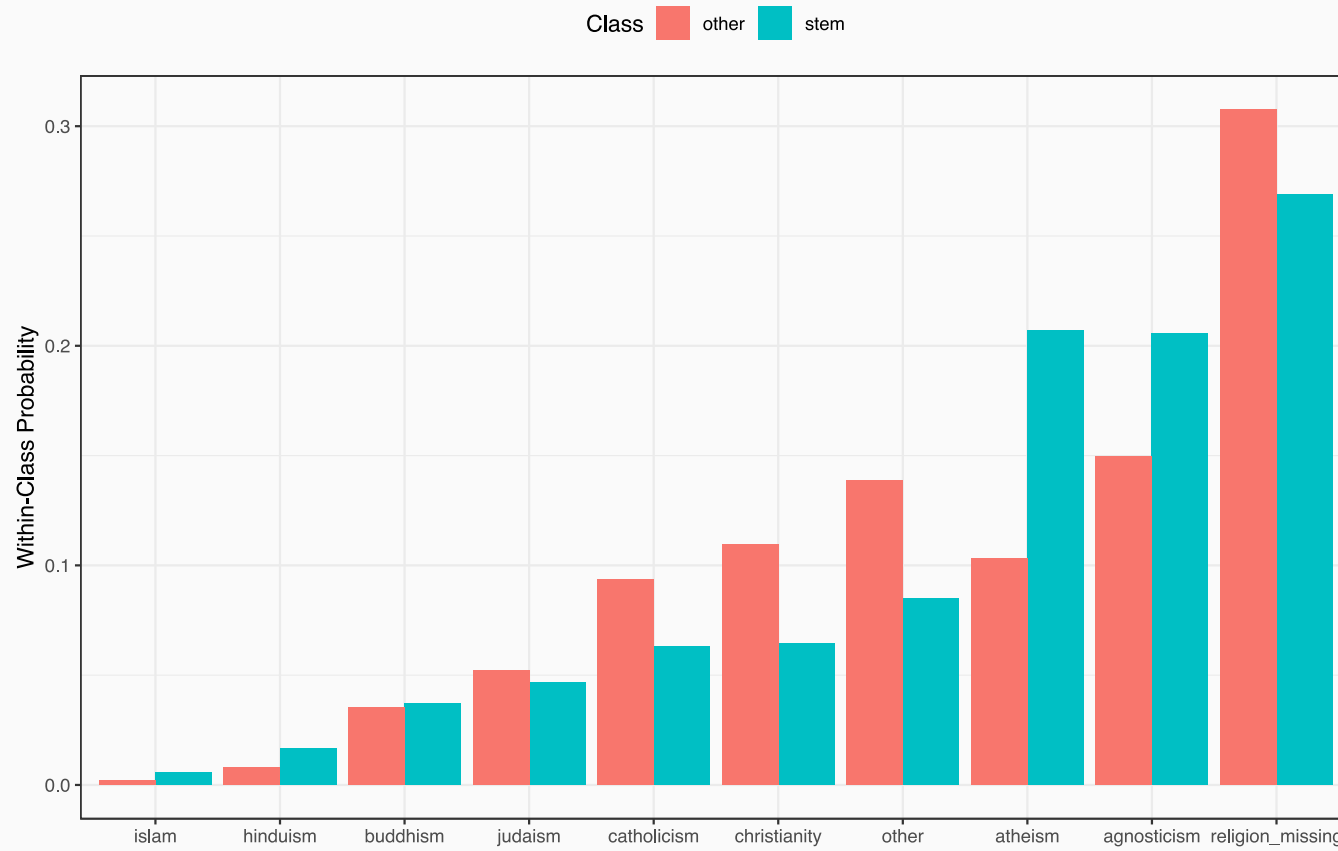
Conditional Densities for Each Class



Conditional Values for Numeric Predictors



Conditional Probabilities for Categorical Predictors



Combining Predictor Scores with the Prior

For an Atheist who wrote 8 words, their likelihood values were:

- $P(\text{Atheist} | \text{8 words}) = 0.518 \times 0.207 = 0.107$
- $P(\text{STEM} | \text{8 words}) = 0.428 \times 0.103 = 0.044$

However, when these are combined with the prior probabilities for each class, the results show:

- $P(\text{Atheist} \text{ and } \text{8 words}) = 0.107 \times 0.182 = 0.02$
- $P(\text{STEM} \text{ and } \text{8 words}) = 0.044 \times 0.818 = 0.036$

We don't need to compute the evidence; we can just normalize these values to add up to 1.

The results is that the probability that this person is in a STEM field is 35.1%.

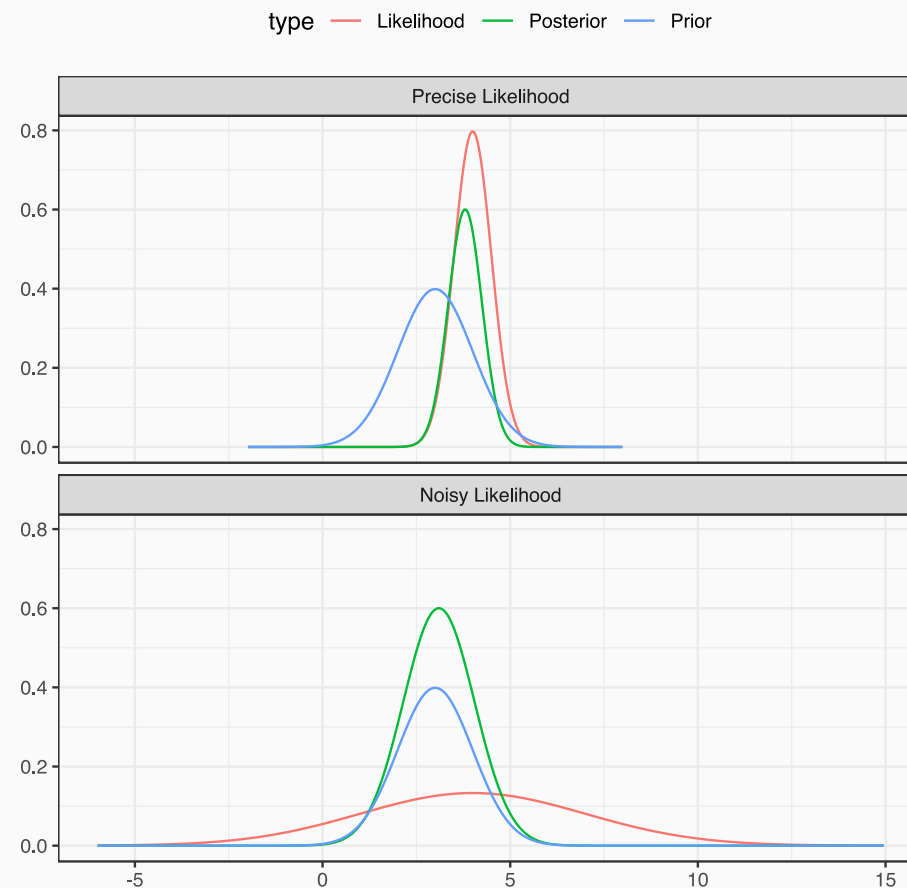
"Shrinkage" Properties of Bayesian Models

When our observed data is of high quality (top panel), the likelihood is narrow and dominates the equation.

However, when our observed information is poor, the likelihood can be very wide and diffuse.

When this occurs, Bayes' rule relies more on our prior belief than on the data in-hand.

This is generally called "shrinkage".



Pros and Cons

Good:

- This model can be very quickly trained (and theoretically in parallel).
- Once trained, the prediction is basically a look-up table (i.e. fast).
- Nonlinear class boundaries can be generated.

Bad:

- Linearly diagonal boundaries can be difficult.
- With many predictors, the class probabilities become poorly calibrated and U-shaped with most values near zero or one.

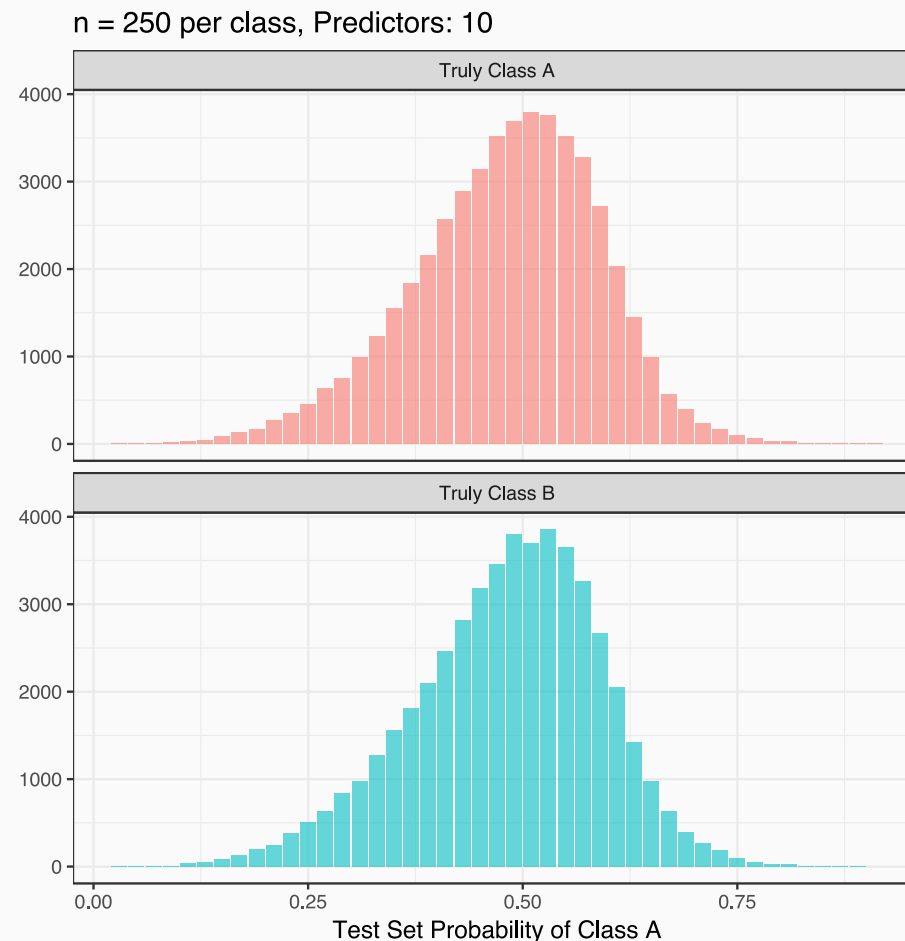
U-Shaped Class Probability Distributions

A completely non-informative data set was simulated using the naive assumption.

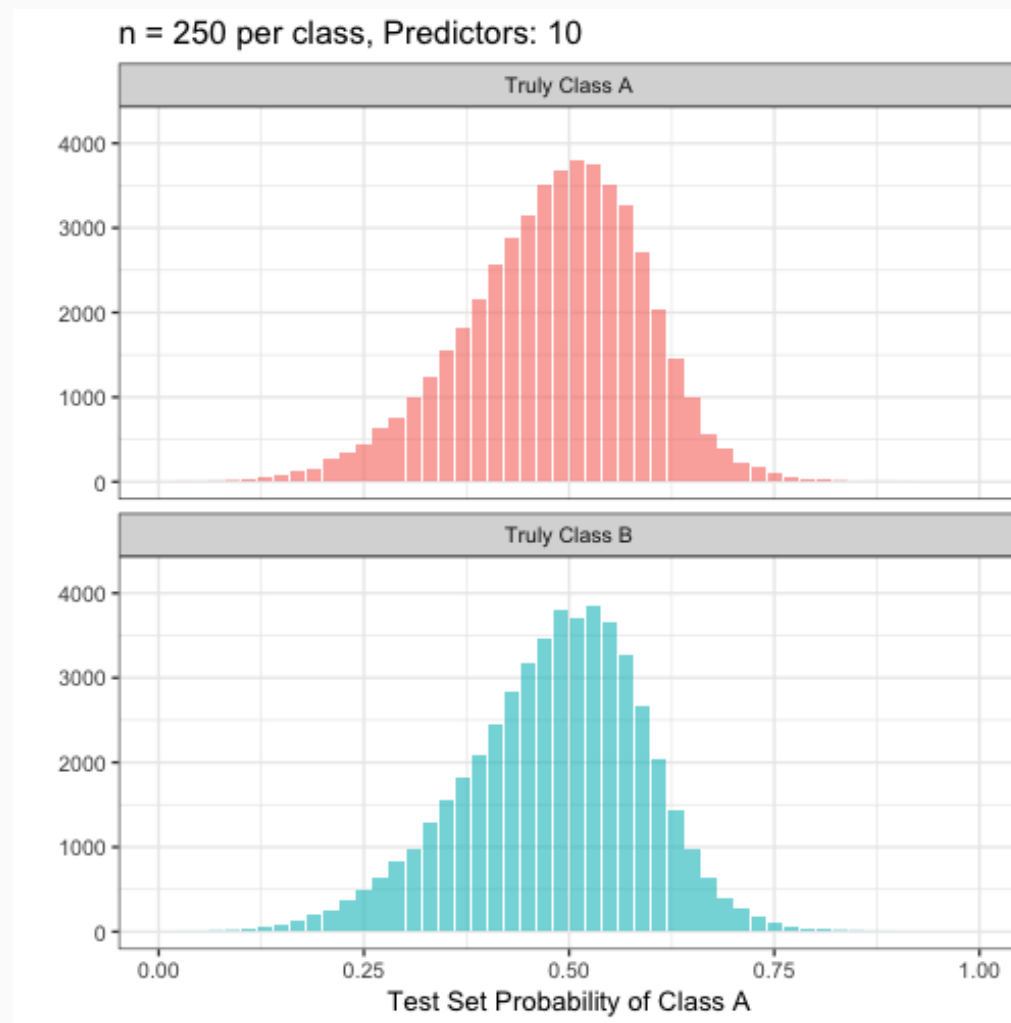
The training set has 500 data points over two classes and 450 predictors.

When a model is fit with 10 predictors, the distribution of the class probabilities gives us shapes that we would expect.

What happens when the number of predictors becomes larger?



U-Shaped Class Probability Distributions



Naive Bayes Data Preparations

There is a step specifically designed for converting binary dummy variables into factors.

First, we identify them, then use quasiquotation (via `quo()`) to pass them to the step function.

```
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
```

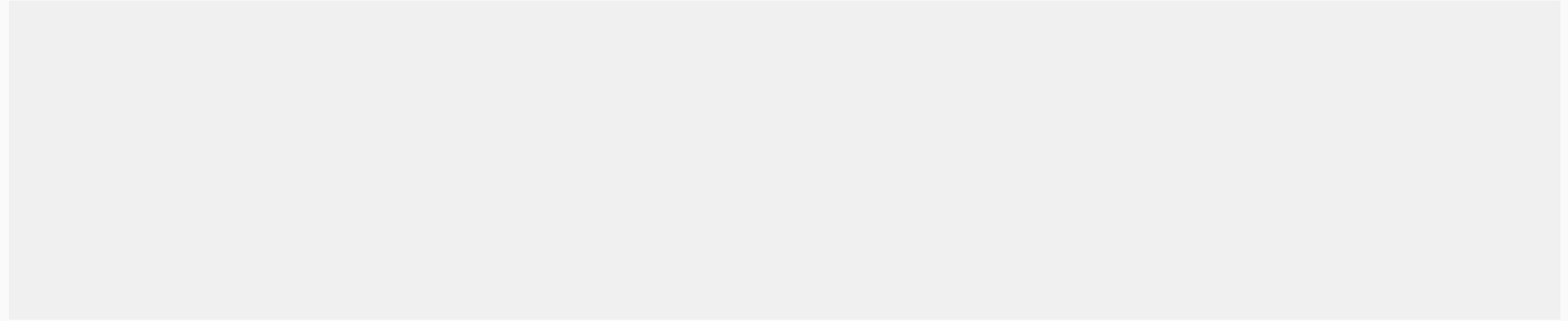
```
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
```

```
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
```

We will use basic defaults for the tuning parameters:

```
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
library(MASS)
```

Naive Bayes Training

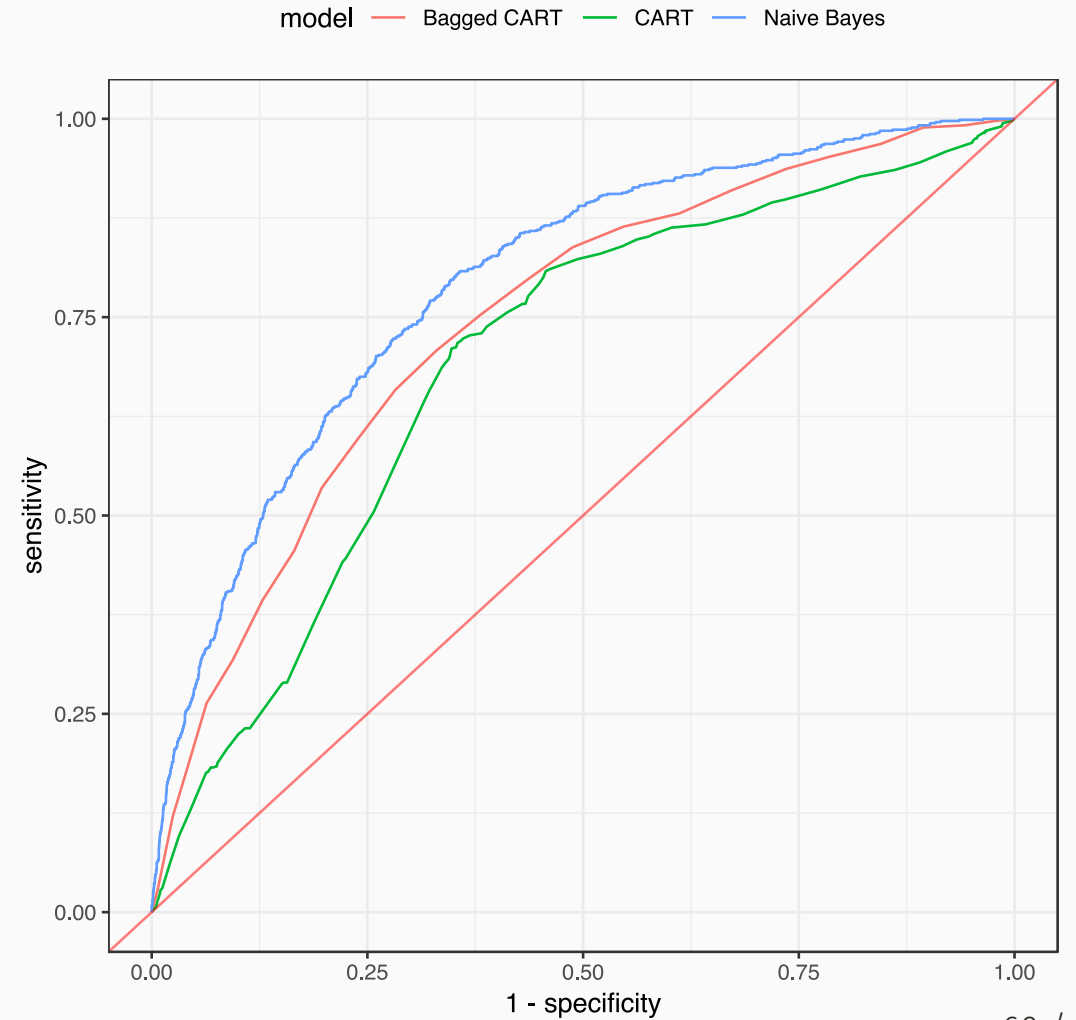
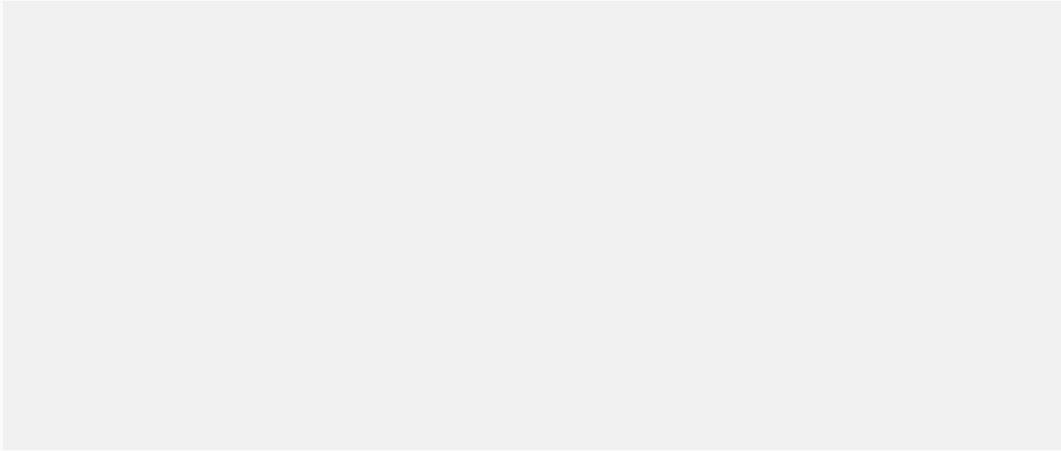


Some warnings:

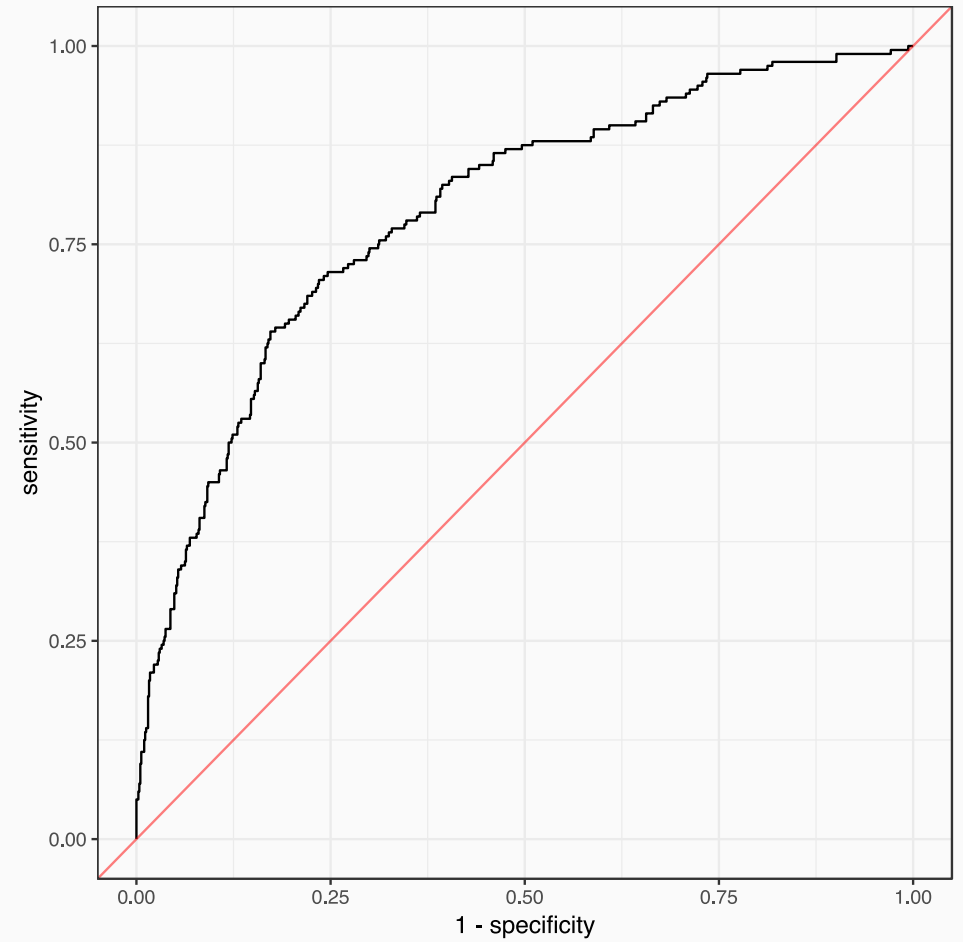
This is due to the poorly calibrated probabilities although the warning is a bit misleading. This issue does not generally affect performance and can be ignored.

Naive Bayes Resampling Profile

Three ROC Curves



Test Set ROC Curve



Test Set Class Probabilities

