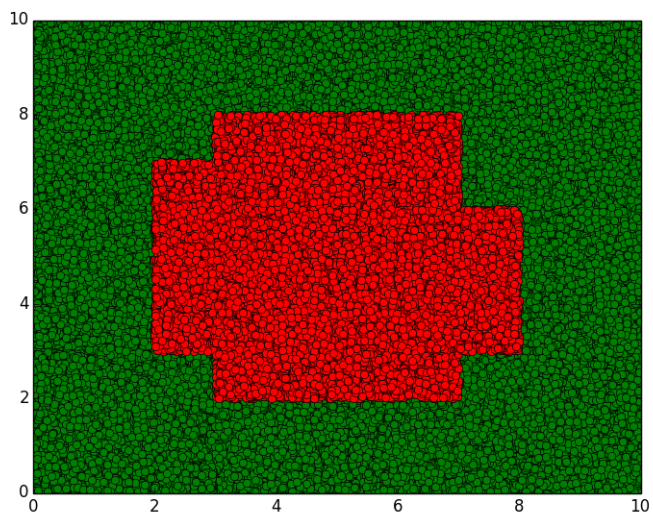Trevor Slaton (tms45)

Joseph Vokt (jpv52)

# 1 Model Selection and Validation

(a) $S_{train}$ accuracy: 22/24. $S_{test1}$ accuracy: 8/10. $S_{test2}$ accuracy: 16/20.

(b) $S_{train}$ accuracy: 24/24. $S_{test1}$ accuracy: 8/10. $S_{test2}$ accuracy: 15/20.

(c) Let $x = \frac{(n_{01}-n_{10})^2}{n_{01}+n_{10}}$ where $n_{01}$ corresponds to the number of cases which were missclassified by the decision tree hypothesis and not the linear hypothesis, and $n_{10}$ corresponds to the number of cases which were missclassified by the linear hypothesis but not the decision tree hypothesis. For $S_{test1}$, $n_{01} = 1$ and $n_{10} = 1$, so $x = 0$. Using matlab chi2cdf(x,1), we get a p-value of 0. Thus we can't say that one model generalizes any better than the other.

(d) For $S_{test2}$, $n_{01} = 3$ and $n_{10} = 2$, so $x = \frac{1}{5} = .2$. Using matlab chi2cdf(x,1), we get a p-value of 0.3453. Thus we can't confidently say that one model generalizes any better than the other.

(e) The accuracy for $S_{test1}$ is the same for both the linear hypothesis and decision tree. The $\chi^2$ statistic for the McNemar test for $S_{test1}$ comparing the two hypotheses confirms that there is no significant difference in how well they generalize. The accuracy for $S_{test2}$ is slightly higher for the linear hypothesis than the decision tree. However, the $\chi^2$ statistic for the McNemar test for $S_{test2}$ comparing the two hypotheses reveals that we are not justified in claiming that one hypothesis generalizes better than the other.
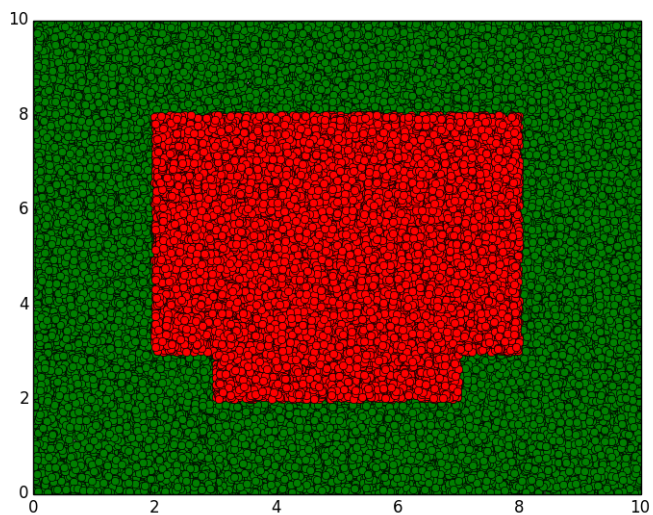
Although the results for the two McNemar tests were different, the conclusion is the same. Unlike for $S_{test1}$, the decision tree hypothesis does separately misclassify more examples in $S_{test2}$ than does the linear hypothesis, but the difference between their misclassification rates is still shown by the McNemar test to be insignificant. In order to determine that results are significant at a confidence level of 95%, the $\chi^2$ statistic for the McNemar test must exceed the critical value 3.84, which mathematically requires both a sizable discrepancy between the separate misclassifications $n_{01}$ and $n_{10}$ and a large total number of misclassifications (e.g. $n_{01} = 40, n_{10} = 60$). In both cases for this question, we have neither.

# 2 Model Averaging with Decision Trees

(a) Individal decision trees:

(b) Averaging decision trees:



(c) Show that a prediction can be expressed as follows:

$$\hat{y} = sign\left(\sum_i^n y_i K(x_{test}, x_i)\right)$$

A test instance, $x_{test}$ can end up in one and only one leaf. Let $k$ be the

number of training instances in that leaf. This means that the similarity measure $K(x_{test}, x_i)$ will be nonzero for exactly $k$ training instances $x_i$. If for example all training instances in that leaf are labeled as $+$, we get

$$\hat{y} = sign\left(\sum_i^k 1/k\right) = sign(1) = +$$

If for example all training instances in that leaf are labeled as $-$, we get

$$\hat{y} = sign\left(\sum_i^k -1/k\right) = sign(-1) = -$$

If the majority of training instances in that leaf are $+$, then the test instance is labeled as $+$. If the majority of training instances in that leaf are $-$, then the test instance is labeled as $-$. If there are an equal number, we will have to flip a coin.

(d) Show that averaging predictions from $M$ different trees,

$$\hat{y} = sign\left(\frac{1}{M}\sum_j^M \sum_i^n y_i K_j(x_{test}, x_i)\right)$$

can be expressed as follows:

$$\hat{y} = sign\left(\sum_i^n y_i \tilde{K}(x_{test}, x_i)\right)$$

Proof:

$$\hat{y} = sign\left(\frac{1}{M}\sum_j^M \sum_i^n y_i K_j(x_{test}, x_i)\right)$$

$$= sign\left(\sum_i^n y_i \frac{1}{M}\sum_j^M K_j(x_{test}, x_i)\right)$$

$$= sign\left(\sum_i^n y_i \tilde{K}(x_{test}, x_i)\right)$$

where

$$\tilde{K}(x_{test}, x_i) = \frac{1}{M}\sum_j^M K_j(x_{test}, x_i)$$

Intuitively, to label test instance $x_{test}$ we are measuring which is greater:
(1) $\sum_i \tilde{K}(x_{test}, x_i)$ where $i$ is such that $x_i$ is labeled with a plus, or
(2) $\sum_i \tilde{K}(x_{test}, x_i)$ where $i$ is such that $x_i$ is labeled with a minus.

$\tilde{K}(x_{test}, x_i)$ represents the average $K$ among all $M$ different trees. If $x_i$ has a $+1$ label, it would contribute to the first term, and if it has a $-1$ label it would contribute to the second. In other words, we label $x_{test}$ based on whether we have more decision trees that predict $+1$ or $-1$.

# 3 Text Categorization with Decision Trees

(a)

(b)

(c)

(d)

(e)

(f)